# Algorithms for understanding the spatial and network organization of biological systems

Uthsav Chitra

March 1, 2024

**PRINCETON**
UNIVERSITY

# Biological systems are organized across a hierarchy of scales: from genes and proteins...


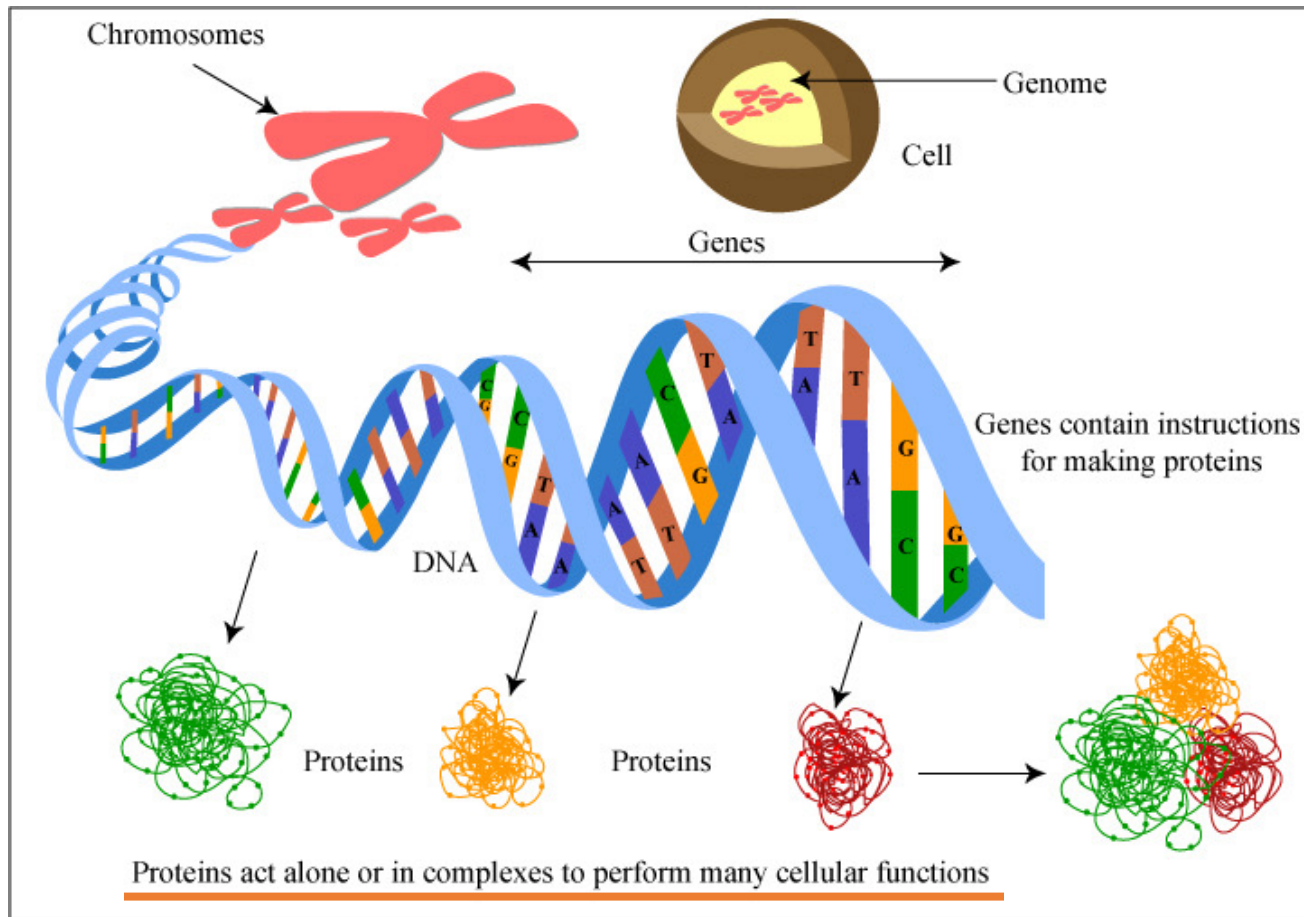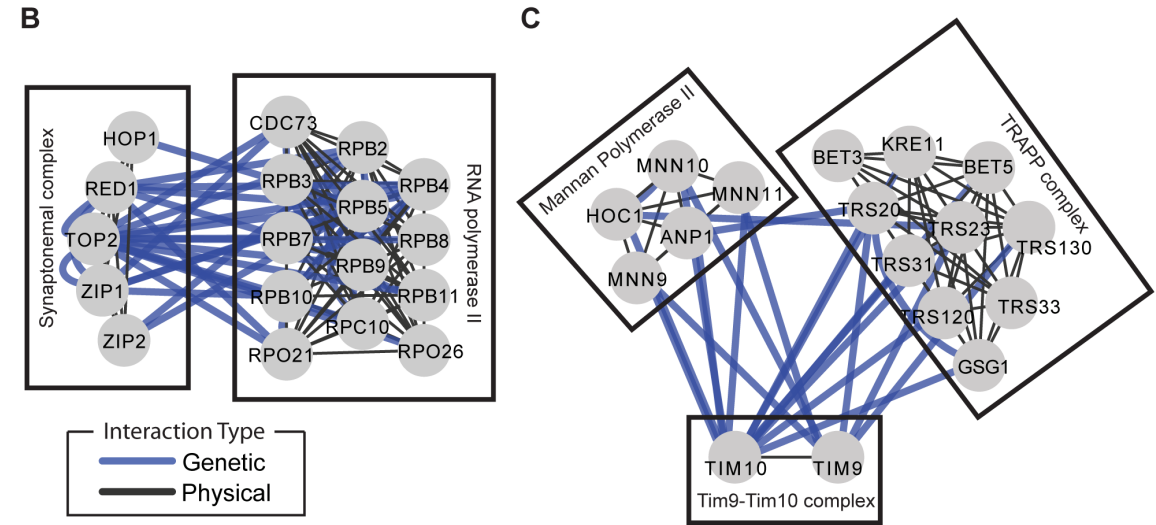
Image from MIT OpenCourseWare
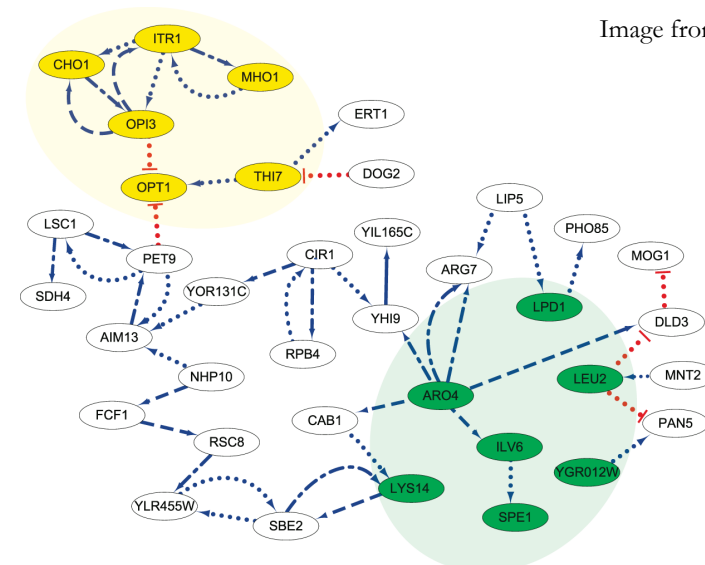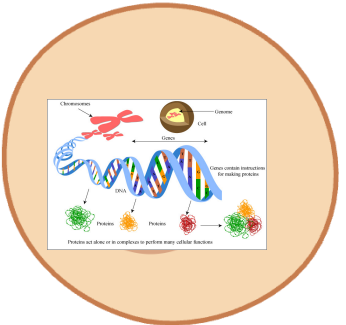


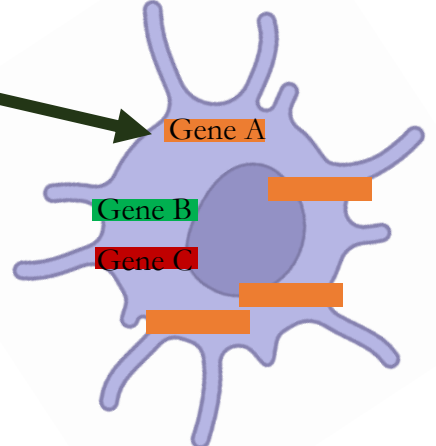Image from Hannum et al., PLOS Genetics 2009



Image from Chen et al, Scientific Reports 2019

# …to cells and tissues

**Fibroblast**
(High expression of Gene A)

**Gene expression**
(e.g. mRNA transcript)

Gene A

Gene B

Gene C

**T-cell**
(High expression of Gene C)

Gene C
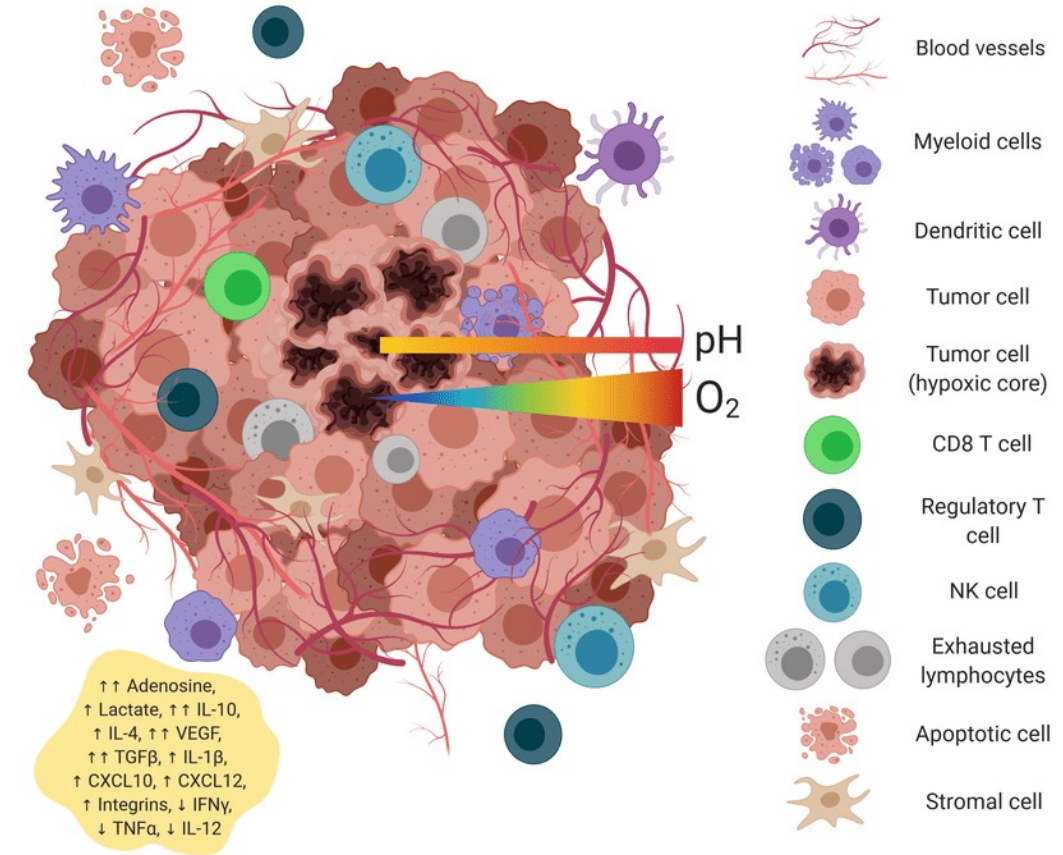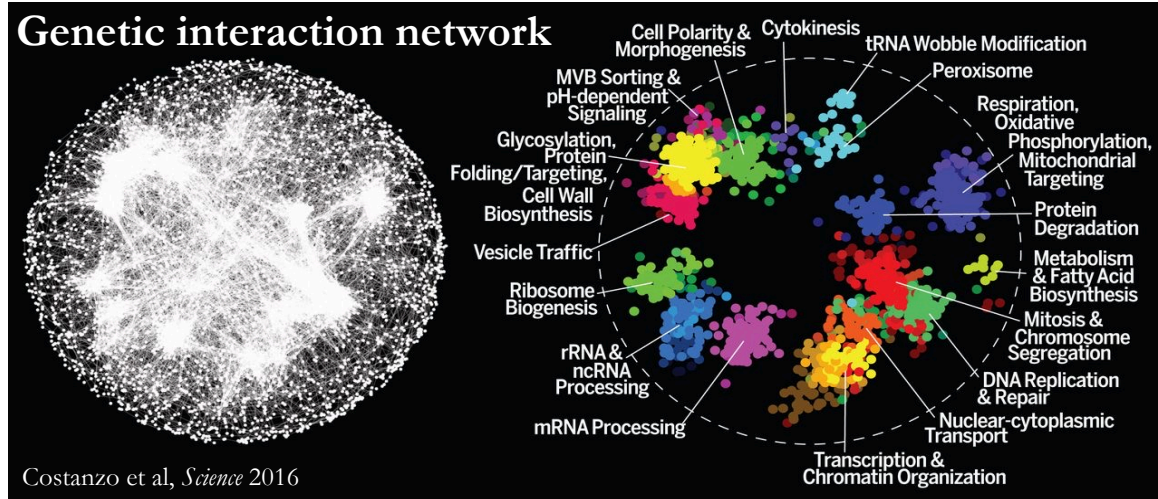
Gene B

Gene A

Spatial arrangement of cells in tissue

# Genetic interactions and cellular organization impact human health and disease



Genetic interaction network

Costanzo et al, *Science* 2016



Cancer mutations alter gene networks

Leiserson et al, Nature Genetics 2015



Fernandez et al, Int J Mol Sci 2019

Spatial heterogeneity in the tumor microenvironment

# High-throughput sequencing data enables study of biological systems



Tissue sample

Sequencing

Genome + mutations
(DNA sequencing)

Gene expression
(RNA sequencing)

Spatial sequencing
(e.g. Slide-Seq, 10x Visium, …)

**Computational methods needed to derive insights from large volume of sequencing data**

# My thesis: computational methods for understanding complex biological systems
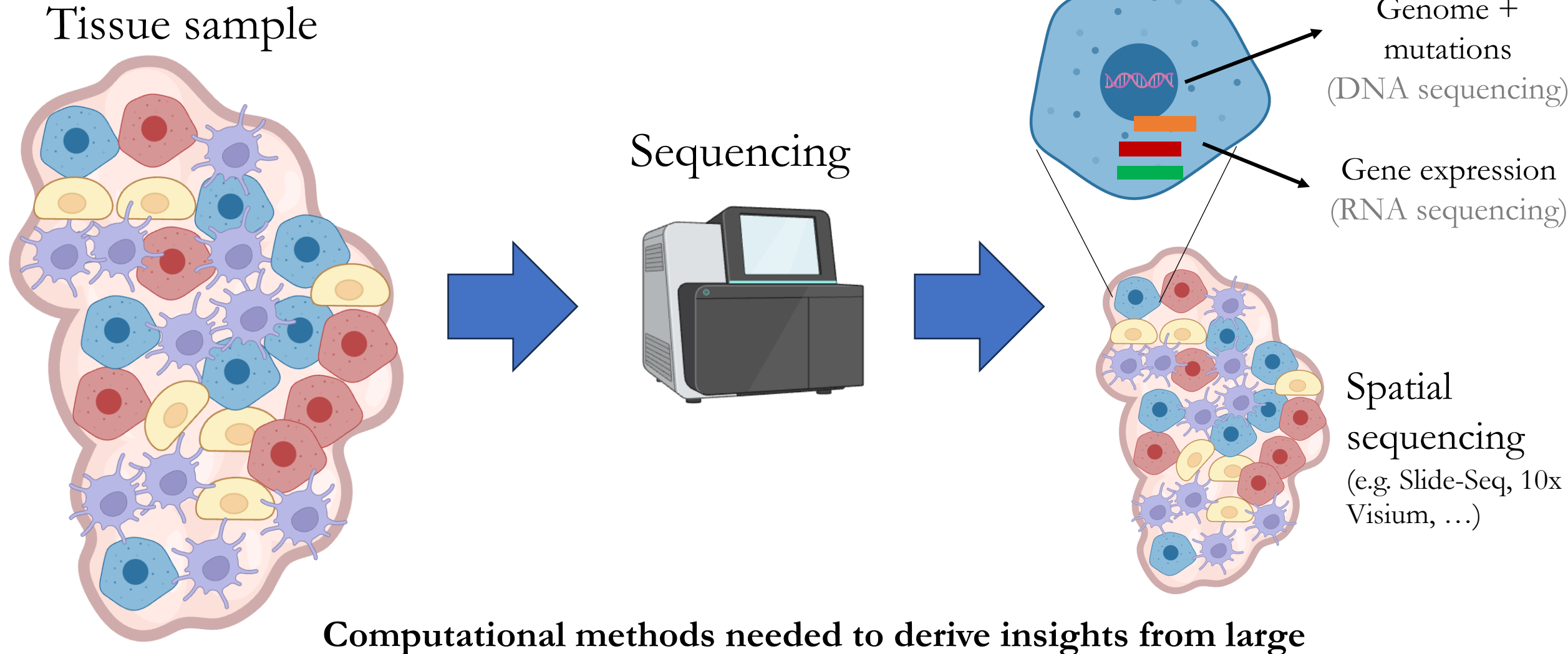
## Spatial biology



Legend:
- Healthy
- Invasive
- Surrounding tumor
- Tumor

## Network interactions and anomalies



NUP53-ASM4

Spatial variation in gene expression
- Ma*, **Chitra***, et al. *RECOMB 2022 + Cell Systems.*
- **Chitra** et al. *RECOMB 2024 + in review at Nature Methods.*

Altered subnetwork identification
- Reyna*, **Chitra***, et al. *RECOMB 2020 + JCB.*
- **Chitra** et al. *ICML 2021.*
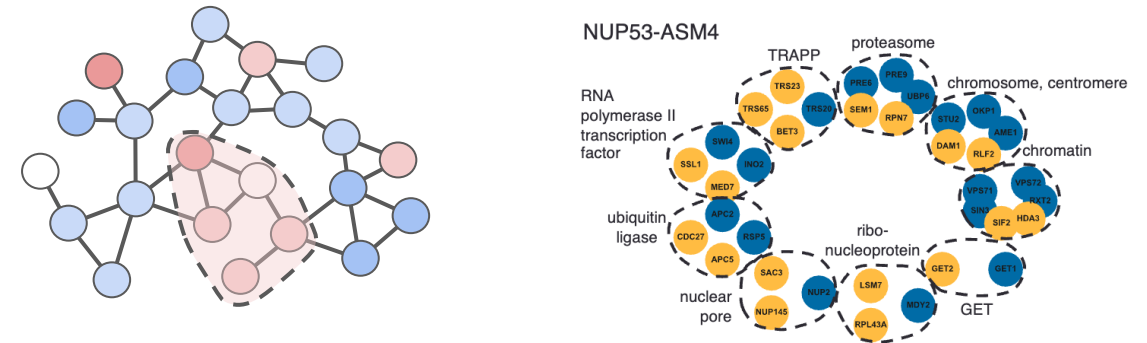- **Chitra***, Park*, Raphael. *RECOMB 2022 + JCB.*

Cell-cell interactions
- Sarkar*, **Chitra***, et al. *In submission at ISMB 2024.*

Learning genetic interactions
- **Chitra***, Arnold*, Raphael. *In review at Nature Genetics.*
- Shuaibi*, **Chitra***, Raphael. *In submission at RECOMB-CCB.*

Machine learning + data mining
- **Chitra** and Raphael. *ICML 2019.*
- **Chitra** and Musco. *WSDM 2020.*

*\* indicates joint first authorship*

# My thesis: computational methods for understanding complex biological systems

## Spatial biology



Healthy
Invasive
Surrounding tumor
Tumor

## Network interactions and anomalies



NUP53-ASM4

### Spatial variation in gene expression
- Ma*, **Chitra***, et al. *RECOMB 2022 + Cell Systems.*
- **Chitra** et al. *RECOMB 2024 + in review at Nature Methods.*

### Cell-cell interactions
- Sarkar*, **Chitra***, et al. *In submission at ISMB 2024.*
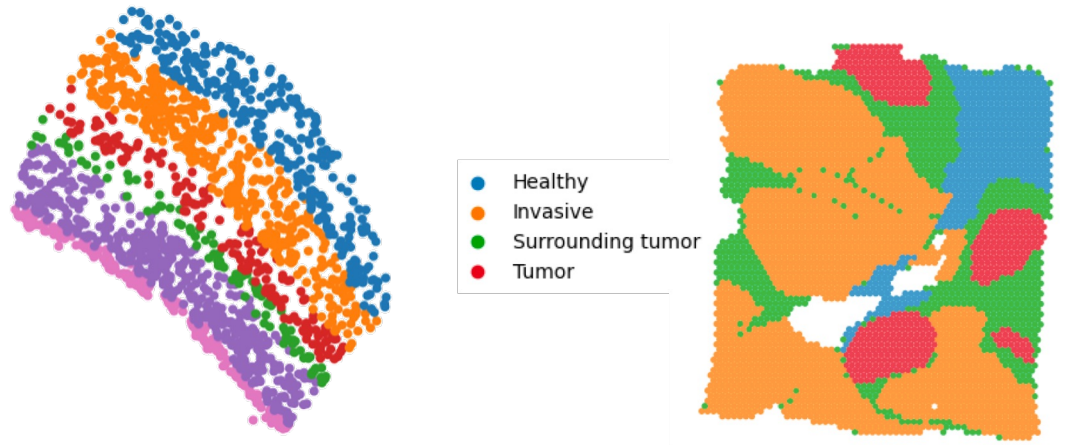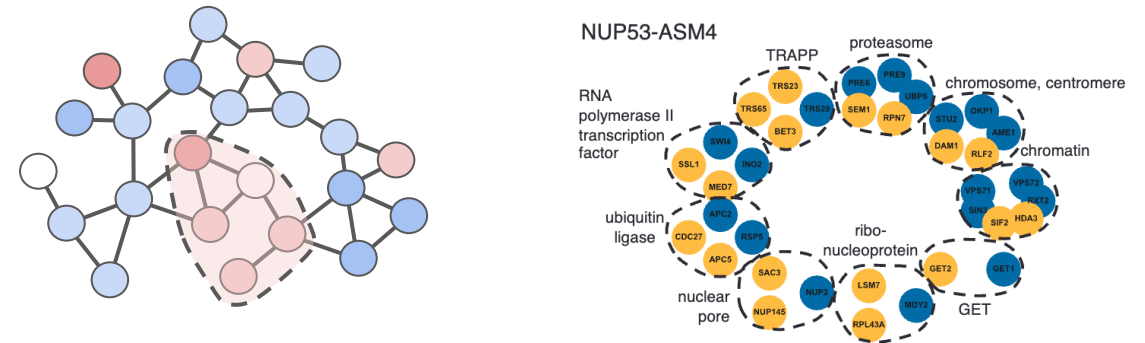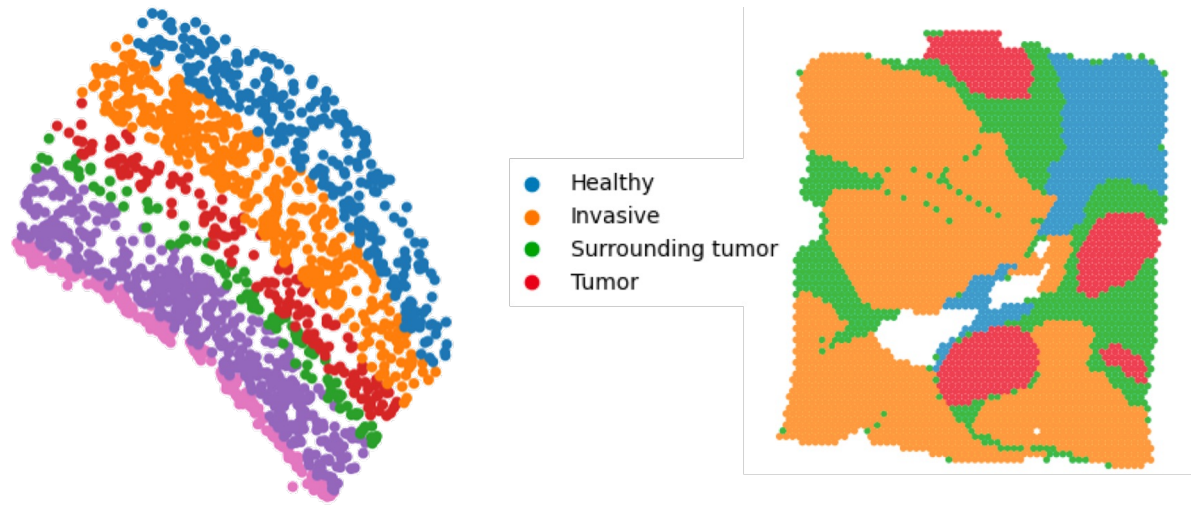
### Altered subnetwork identification
- Reyna*, **Chitra***, et al. *RECOMB 2020 + JCB.*
- **Chitra** et al. *ICML 2021.*
- **Chitra***, Park*, Raphael. *RECOMB 2022 + JCB.*

### Learning genetic interactions
- **Chitra***, Arnold*, Raphael. *In review at Nature Genetics.*
- Shuaibi*, **Chitra***, Raphael. *In submission at RECOMB-CCB.*

### Machine learning + data mining
- **Chitra** and Raphael. *ICML 2019.*
- **Chitra** and Musco. *WSDM 2020.*

* indicates joint first authorship   : new work since pre-FPO

# Modeling spatial variation in gene expression



Ma*, **Chitra*,** Zhang, Raphael. *RECOMB 2022 + Cell Systems.*

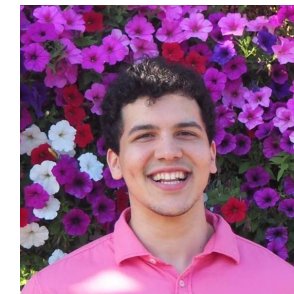**Chitra** et al. *RECOMB 2024 + in review at Nature Methods.*

* indicates joint first authorship

Cong Ma

Shirley Zhang

Brian Arnold

Hirak Sarkar

Sereno Lopez-Darwin

Kohei Sanno

Ben Raphael

# Spatially Resolved Transcriptomics (SRT)



Tissue sample

Barcoded Grid of Spots

6.2 mm

6.6 mm

1007 spots

Berglund et al. Nat Com. 2018

●●●●● |100 μm
200 μm

RNA sequencing

**Gene expression matrix + spatial coordinates**

Genes (G)

Spatial (2D)

Spot (S)

$\lambda_s$

**Technologies**: Slide-Seq, 10x Visium, MERFISH, …

**High-throughput:** measure 1,000-20,000 genes at 1,000-10,000 spatial locations (each spot contains 1-20 cells)
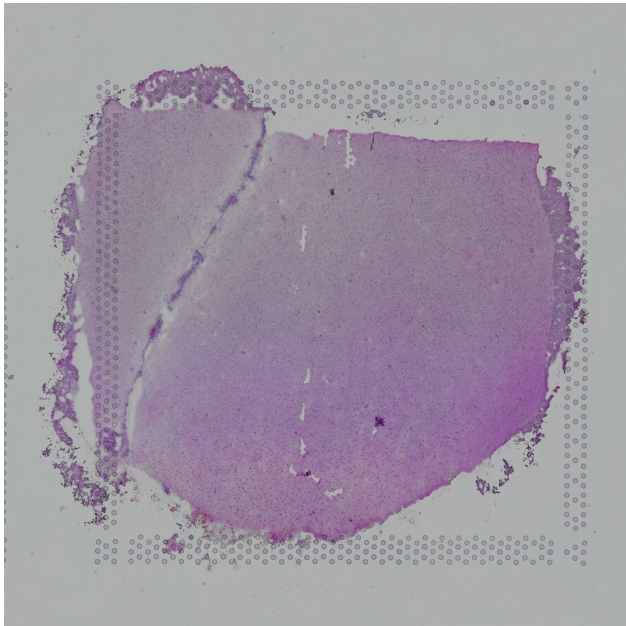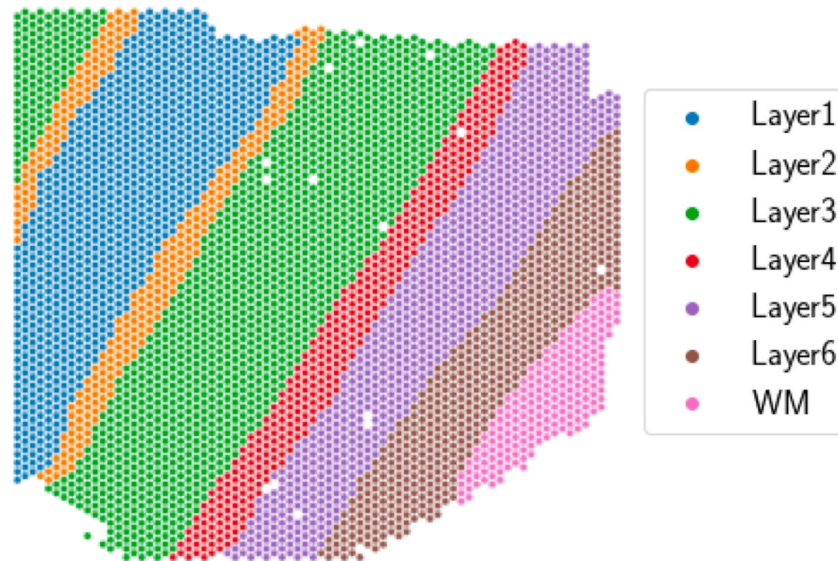
# Spatially resolved transcriptomics (SRT) reveals new biology

## Spatial domains/cell types

H&E stain

Neuronal tissue layers

- Layer1
- Layer2
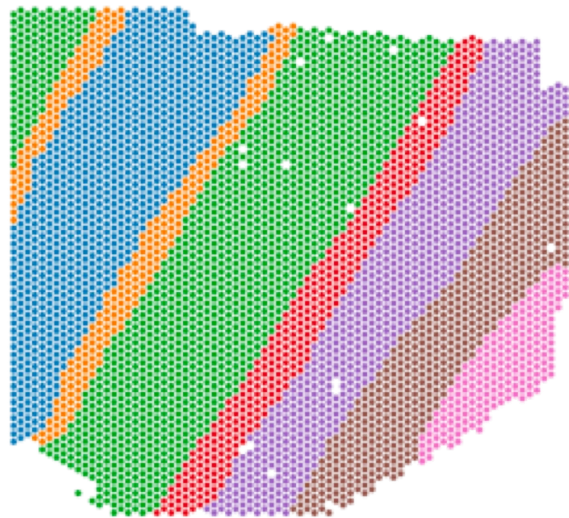- Layer3
- Layer4
- Layer5
- Layer6
- WM

## Marker genes

(differentially expressed across domains)

Layer5 marker gene *PCP4*

UMI count

Human dorsolateral pre-frontal cortex (DLPFC) [Maynard et al., *Nat Neurosci* 2021]

# Challenge: SRT data is **very sparse**!



Sample 151508

Layer1
Layer2
Layer3
Layer4
Layer5
Layer6
WM

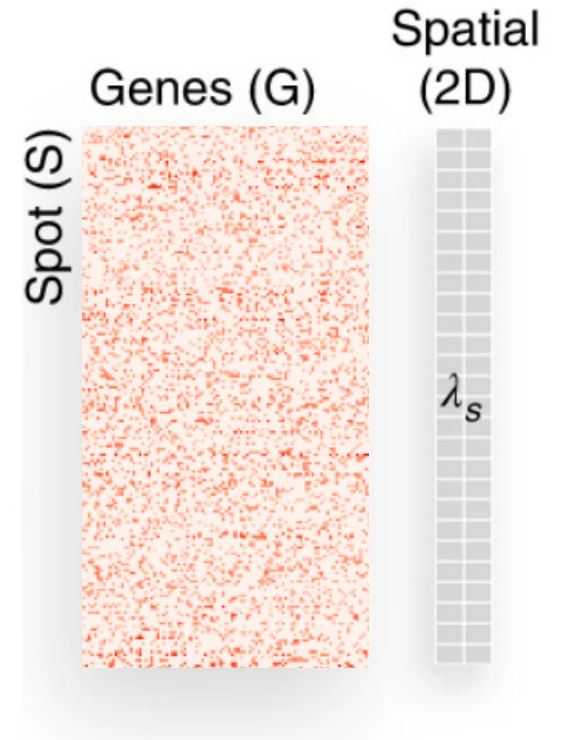Sparse expression of marker gene *TRABD2A*

UMI count

Median gene has non-zero
expression in <5% spots
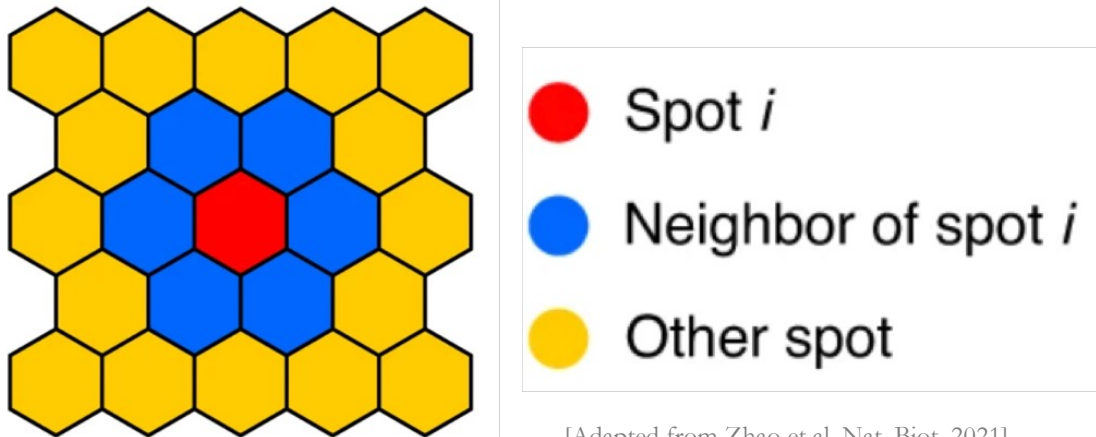
Genes (G)

Spatial (2D)

Spot (S)

$\lambda_s$

**Sparse** matrix:
>90% zeros

# Overcoming sparsity by incorporating spatial information

Most algorithms use **local models**: nearby spots
have similar cell type / expression

- **Hidden Markov Random Field (HMRF)**: BayesSpace [Nat. Biotech 2021],
  SPICEMIX [Nature Genetics 2022], Giotto [Genome Biology 2021], scGCO [Nat Comm 2022] …

- **Graph neural networks (GNN)**: SpaGCN [Nature Genetics 2021], STAGATE
  [Nature Communications 2022], SEDR [Genome Med 2024] ,…

- **Gaussian Processes**: SpatialDE [Nature Methods 2018], SPARK [Nature Methods 2020],
  SPARK-X [Genome Biology 2021] , nnSVG [Nat Comm 2023] …



**Hidden Markov
Random Fields**



Yuan and Bar-Joseph,
Genome Biol. 2020

**Graph Neural
Networks**



Svensson *et al.*, Nat. Methods 2018

**Gaussian Processes**



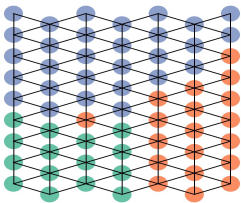[Adapted from Zhao et al. Nat. Biot. 2021]

# Overcoming sparsity by incorporating spatial information

Most algorithms use **local models**: nearby spots have similar cell type / expression
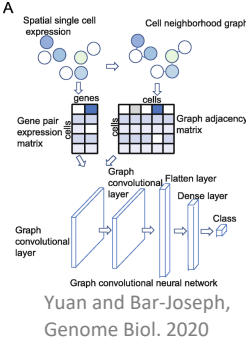
- **Hidden Markov Random Field (HMRF)**: BayesSpace [Nat. Biotech 2021], SPICEMIX [Nature Genetics 2022], Giotto [Genome Biology 2021], scGCO [Nat Comm 2022] …

- **Graph neural networks (GNN)**: SpaGCN [Nature Genetics 2021], STAGATE [Nature Communications 2022], SEDR [Genome Med 2024] ,…

- **Gaussian Processes**: SpatialDE [Nature Methods 2018], SPARK [Nature Methods 2020], SPARK-X Genome Biology 2021] , nnSVG [Nat Comm 2023] …

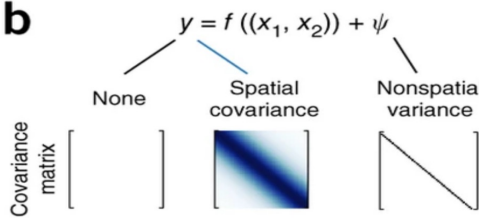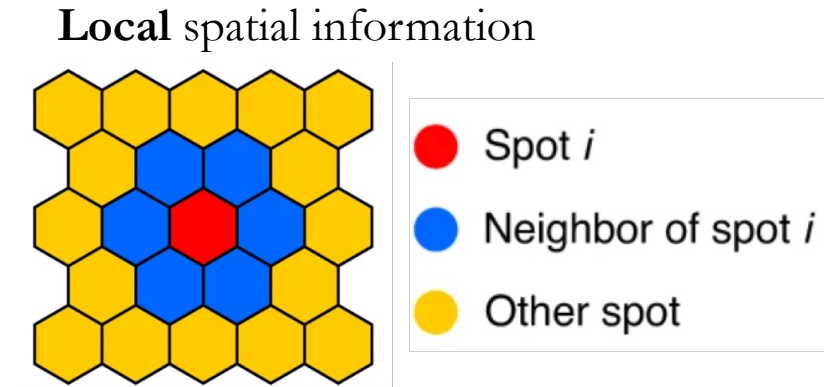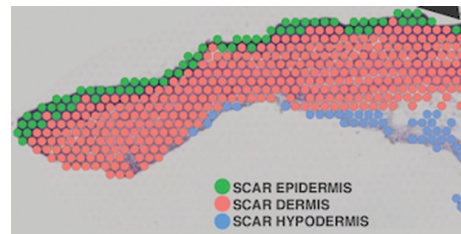**Local** spatial information

- Spot *i* (red)
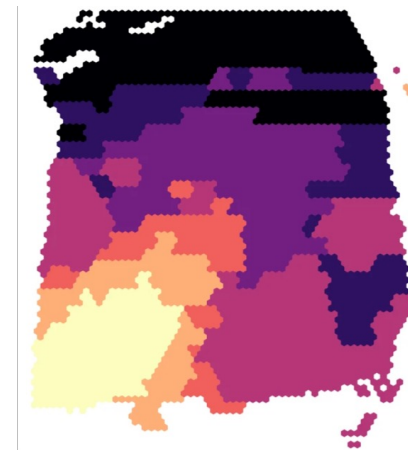- Neighbor of spot *i* (blue)
- Other spot (yellow)

[Adapted from Zhao et al. Nat. Biot. 2021]
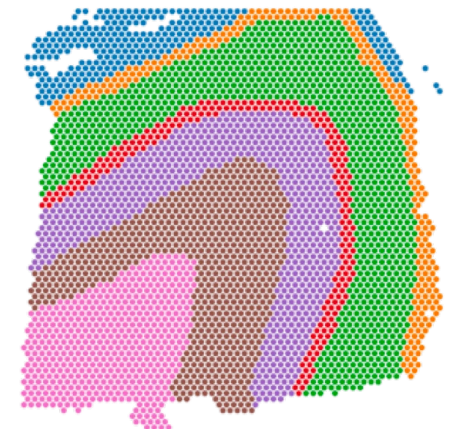
## **Global model:** Cortex is made of layers!

- Many layered tissues: skin, ureter, eye, …
- Can we incorporate the <u>layered geometry</u> in a gene expression model?

Skin
(Foster et al., 2021)
SCAR EPIDERMIS
SCAR DERMIS
SCAR HYPODERMIS

Spatial domains learned
by Giotto (local model)

Annotated layers

# A simple layered tissue

DLPFC sample 151508
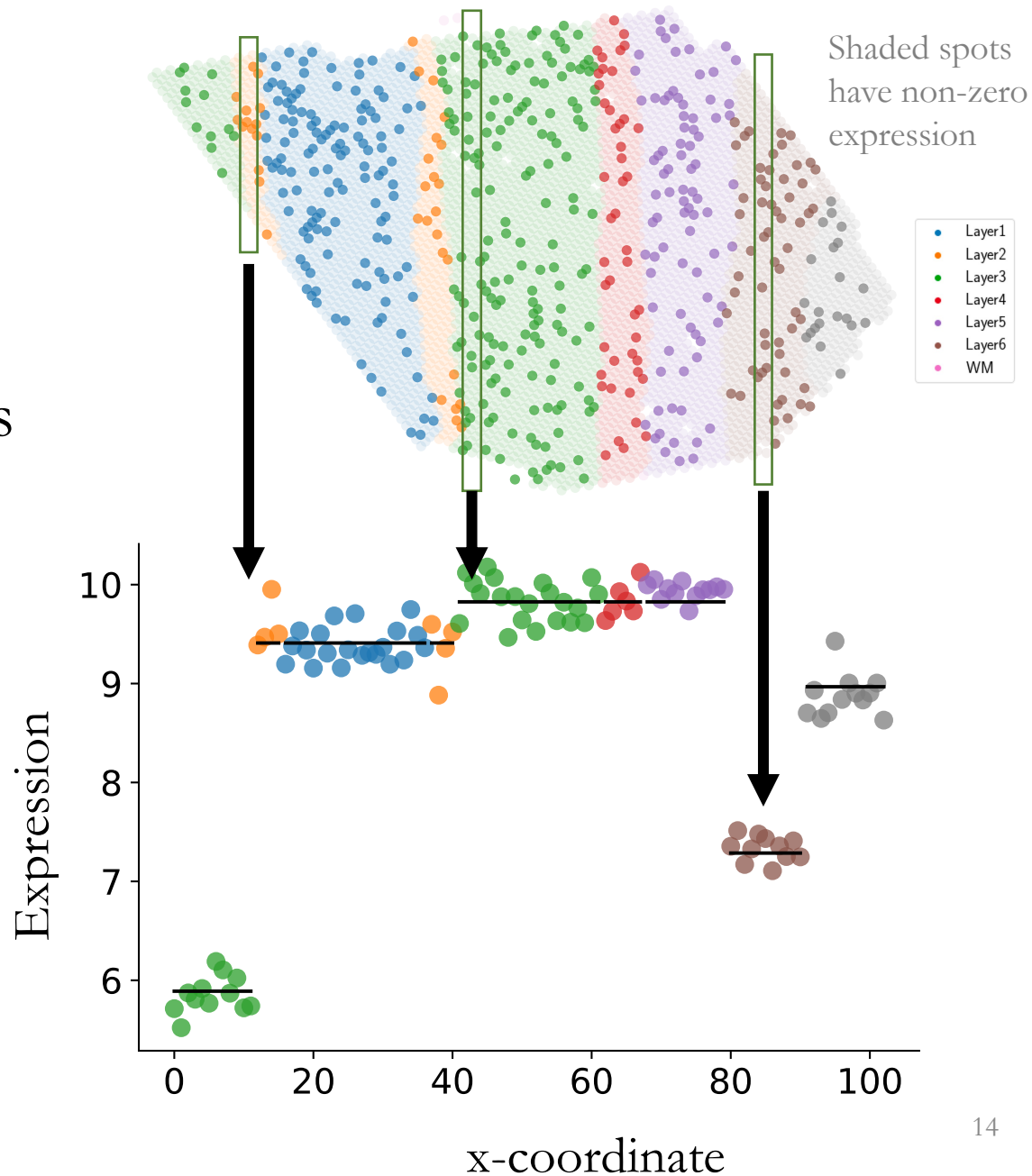(approximately axis aligned)



$f_g(x, y) = f_g(x)$ is piecewise constant

(Marker) gene expression is
$\approx$***constant*** along $y$-axis

expression only depends on x-coord = distance to layer boundary (layer depth)

Pool sparse expression along y-axis



Shaded spots have non-zero expression

Expression

x-coordinate

14

# A simple layered tissue

DLPFC sample 151508
(approximately axis aligned)

Shaded spots
have non-zero
expression

$f_g(x, y) = f_g(x)$ is
**piecewise linear**

Layer1
Layer2
Layer3
Layer4
Layer5
Layer6
WM

$y$

$x$

(Marker) gene expression is
$\approx$***constant*** along $y$-axis

expression only depends on x-coord = distance to layer boundary (layer depth)

Expression

−9.0

−9.2

−9.4

−9.6

x-

15

# How to model layer depth in tissues with complex layered geometry?



**Global layered tissue geometry**

**Expression function:** Piecewise linear function of <u>layer depth</u>

# Conformal maps and harmonic functions model layer depth



A **conformal map** $\Phi: D \subseteq \mathbb{C} \to \mathbb{C}$ **locally preserves angles** between curves

$\mathrm{u} = \mathrm{Re}\ \Phi$ is a **harmonic function** (satisfies heat eq)

Layer depth = **isotherms** (contours) of heat equation

# Layered tissue problem formulation



## Input

- Spot coordinate $\mathbf{s}_i = (x_i, y_i)$.
- Transcript count matrix $A = [a_{i,g}]$, $a_{i,g}$ for $i^{th}$ spot and $g^{th}$ gene.
- Number of layers $L$.

## Probabilistic model

$$a_{i,g} \sim \text{Poisson}\left( C_i \exp(f_g(\mathbf{\Phi}(x_i, y_i))) \right)$$

Piecewise linear expression: $f_g$

Global layered geometry: layer boundaries at $b_1, b_2, \ldots, b_{L-1}$

## Maximum likelihood objective

$$\underset{\substack{\text{breakpoints } b_1 < b_2 < \cdots < b_{L-1} \\ \text{piecewise linear } f_1, \ldots, f_G \\ \text{conformal maps } \mathbf{\Phi} = (\Phi_1, \ldots, \Phi_L)}}{\arg\max} \sum_{g=1}^{G} \left( \sum_{i=1}^{N} \log P\left(a_{i,g} \mid f_g(\mathbf{\Phi}(x_i, y_i))\right) \right)$$

# Solutions to special cases

**Case 1:** Approximate layer boundaries $\widetilde{\Gamma}_i$ are given

Step 1: Construct conformal map(s) $\Phi$ by solving heat equation



Step 2: Segmented regression

[Bai and Perron, *Econometrica* 1998]

- dynamic programming algorithm in O(LN²G) time
- L=#layers, N=#spots, G=#genes



**Case 2:** Layer boundaries $\Gamma_i$ are non-intersecting lines (not given)

DP algorithm (~Nussinov algorithm for RNA folding) to find best lines $\Gamma_i$, conformal maps $\mathbf{\Phi}$, piecewise functions $f_g$

# Overview of Belayer

# Belayer accurately identifies cortical layers in human DLPFC



*Takeaway: Belayer (**global** model) outperforms **local** models*

# Belayer identifies marker genes in human DLPFC data

We identify marker genes using slopes/discontinuities of gene expression functions $f_g$



**Belayer**



**(A)**

**(B)**

SpatialDE (AUPRC=0.017)

SPARK (AUPRC=0.029)

HotSpot (AUPRC=0.029)

**Belayer:** Max slope (AUPRC=0.032)

**Belayer:** Layer 3 slope (AUPRC=0.116)

Known marker gene. Involved in neuronal damage.

Not a recorded marker gene. Involved in neuronal signaling.

*Takeaway: Belayer (**global** model) outperforms **local** models*

# Mouse skin wound (10x Visium)



Approximate layer boundaries

Layer depth

Manual annotation

Belayer ARI=0.52

E
D
H

Retnla expression — Layer depth

Trdn expression — Layer depth

Thbs2 expression — Layer depth

Spatially varying genes involved in muscle contraction and wound healing

# Summary of Belayer

**Global** model of spatial gene expression for layered tissues
piecewise linear functions + conformal maps

Belayer simultaneously learns

- **tissue geometry** (layers) and
- **spatially varying genes** (slopes of piecewise linear functions)

from sparse SRT data



Ma*, Chitra*, et al. *Cell Systems* 2022 [Also: RECOMB 2022]

Paper          Code

25

# Limitations of Belayer



**Belayer: Geometric Model**

$$u = \mathbf{Re}\,\Phi$$
$$\Delta u = 0$$

$u(\Gamma_2) = 0$

$u(\Gamma_3) = 1$

Conformal maps
Harmonic functions
Heat equation
…

0   1

✅ **Special case 1:** Manually annotated layer boundaries

✅ **Special case 2:** Layer boundaries are lines

❌ **Cannot model**: tissues with complex or unknown boundaries

Belayer is **<u>supervised</u>** – requires *prior knowledge* of tissue geometry.

Mouse cerebellum

# A new approach



**Belayer: Geometric Model**

$u(\Gamma_2) = 0$

$u(\Gamma_3) = 1$

$\xrightarrow{\substack{u = \mathbf{Re}\, \Phi \\ \Delta u = 0}}$

Conformal maps
Harmonic functions
Heat equation
…

0    1

**Deep learning**

$x$

$y$

$\phi$

$a_1$
$a_2$
$a_3$
$a_4$
$\cdots$
$a_G$

✅ **Special case 1:** Manually annotated layer boundaries

✅ **Special case 2:** Layer boundaries are lines

❓ learn layer depth without prior knowledge

Mouse cerebellum

# **GASTON**: neural network architecture



Coordinates

Gene expression

**GASTON**

Gradient Analysis of Spatial Transcriptomics
Organization with Neural networks

**Topographic map** of tissue slice

Isodepth

**Input:** Spatial coordinates

$x$ $y$

$\phi$

**Hidden layer:** isodepth

$a_1$ $a_2$ $a_3$ $a_4$ $\cdots$ $a_G$

**Output:** Gene expression

Training is **unsupervised**!
(like an auto-encoder)

# **Isodepth** defines "topography" of gene expression

Coordinates

Gene expression



GASTON

Gradient Analysis of Spatial Transcriptomics Organization with Neural networks

**Topographic map** of tissue slice



Isodepth

Training is **unsupervised**! (like an auto-encoder)

**Input:** Spatial coordinates

$x$ $y$

**Hidden layer:** underline{isodepth}

$\phi$

$a_1$ $a_2$ $a_3$ $a_4$ $\cdots$ $a_G$

**Output:** gene expression



**Isodepth** = contours of equal depth: $\phi = c$

- Generalizes *relative depth* from Belayer

- *Neural field* model (used in computer vision/graphics)

Spatial gradients $\nabla \phi$ (gradient of isodepth)

- Directions of maximum change in gene expression

- Gradient field $\nabla \phi$ is *"conservative"* (no curl)

Gene expression functions $f_g(\phi(x, y))$

# Human DLPFC: GASTON outperforms other neural networks and unsupervised Belayer

**GASTON isodepth**



0.99 correlation with (supervised) Belayer depth!



**Manual annotation**

**GASTON**

**Belayer (supervised)**

**Belayer (unsupervised)**

- ■ DLPFC/GASTON/Belayer Layer 1
- ■ DLPFC Layer 2
- ■ DLPFC/GASTON/Belayer Layer 3
- ■ DLPFC Layer 4
- ■ DLPFC/GASTON/Belayer Layer 5
- ■ DLPFC/GASTON/Belayer Layer 6
- ■ DLPFC/GASTON/Belayer White Matter (WM)

33

# GASTON: Mouse Cerebellum (Slide-seqV2)

**Topographic map**

**Spatial domains**



Genes (G)

Spatial (2D)

Spot (S)

$y_{sg}$

$\lambda_s$

**GASTON**

23,096 genes ✕ 9,985 spots

Cable et al., *Nature Methods* 2022

Neuronal layers/cell types
- Oligodendrocytes
- Granule
- Purkinje
- Bergmann
- MLI1
- MLI2
- Other cell types

Purkinje-Bergmann layer

Molecular layer

e

Oligodendrocyte layer

Granular layer

Purkinje–Bergmann layer

Molecular layer

Baizer, Front. Hum. Neurosci. 2014

Median spot has 370 UMIs

*SBK1* expression

Median gene is expressed in 0.2% of spots

Oligodendrocyte layer

Granule layer

P-B layer

Molecular layer

Expression

Isodepth

GASTON *SBK1* expression function

Spatial coherence score

GASTON    NSF    SpaGCN    SpiceMix

36

# Cell type and gene expression gradients



**GASTON**

Marker gene identification

**Cell type-attributable gradient**

*CALB1*

**Other-attributable gradient**

*CAMK2B*

# Olfactory bulb (Stereo-seq)

9,825 spots ✕ 27,106 genes

- Olfactory nerve layer (ONL)
- Glomerular layer (GL)
- External plexiform layer (EPL)
- Mitral cell layer (MCL)
- Internal plexiform layer (IPL)
- Granule cell layer (GCL)
- Rostral migratory stream (RMS)

**DAPI**

**Isodepth and (negative) spatial gradients**

**SpaceFlow** (diffusion pseudotime)

**GASTON**

**SpaGCN**



**Cell type-*attributable* gradient**

**Other attributable gradient**



38

# GASTON identifies gradients in tumor microenvironment

Colorectal tumor slice (stage IV)
(Wu et al, Cancer Discovery 2022)



**GASTON: spatial domains + isodepth**



High isodepth

Low isodepth

- Domain 1 (tumor)
- Domain 2 (tumor-adjacent stroma)
- Domain 3
- Domain 4
- Domain 5

Intrastromal variation

Discontinuity

Intratumoral variation



564 | 56 | 93
50 | 67
186
556

- Stroma
- Tumor

**Intratumoral variation**
aerobic metabolism



*COX7B*

**Discontinuity/Intrastromal variation**-
epithelial-mesenchymal transition (EMT)



*ACTA2*

*TAGLN*

# Summary: **GASTON**

- Isodepth describes **topographic map** and **spatial gradients** of gene expression within tissue slice

- GASTON: **unsupervised** deep learning algorithm to learn isodepth
  - Uncovers spatial domains and gradients of gene expression/cell type

Chitra et al. In review at *Nature Methods*
[Also: RECOMB 2024]

Preprint          Code



Coordinates

Gene expression

GASTON
Gradient Analysis of Spatial Transcriptomics
Organization with Neural networks

**Topographic map** of tissue slice

Isodepth

# Identifying altered subnetworks (network anomalies)



Matt Reyna

Rebecca Elyanow

Tyler Park

Kimberly Ding

Jasper C. H. Lee

Ben Raphael

Reyna*, **Chitra***, Elyanow, Raphael. *RECOMB 2020 + Journal of Computational Biology.*

**Chitra**, Ding, Lee, Raphael. *ICML 2021.*

**Chitra***, Park*, Raphael. *RECOMB 2022 + Journal of Computational Biology.*

* indicates joint first authorship

# Spatial anomaly detection

Spatial anomaly in biology:

Spatial anomaly in epidemiology:



Tumor detection (prostate cancer)

Disease hotspots identification (breast cancer incidence, NYC)

# **Network** anomaly detection



6.2 mm

6.6 mm

1007 spots

100 μm
200 μm

Cancer

Expected % of reads
50
40
30
20
10
0

**Tumor detection**
**(prostate cancer)**

Epidemiology:
**Disease hotspots**
**(breast cancer, NYC)**

High score

Low score

## **Network anomaly detection**

- Vertices = spatial locations (e.g. tissue spots, census tracts)
  - Edges connect spatially adjacent vertices
- Score: % cancer cells OR disease incidence OR …

Anomalies: subnetworks with large score

# Protein-Protein Interaction (PPI) Networks

Vertices: proteins
Edges: physical interactions between proteins

# Altered Subnetwork Problem (ASP)

(also called <u>network anomalies</u>, <u>network modules</u>, <u>active subnetworks</u>,)



High

Low

**Given**:

1) Interaction network G = (V,E)
2) Vertex scores $X_v$
   - eg from mutations, differential expression, …

**Goal**: Identify **high-scoring subnetworks** of G ("altered subnetworks")

# Altered subnetworks reveal interacting genes relevant to complex traits+diseases



Somatic mutations in cancer

Leiserson, Vandin et al (Nature Genetics 2015)

Complex traits (e.g. height, diabetes, …)

46

# Many algorithms developed over past 20 years for identifying altered subnetworks

**Table 1 | Some recent bioinformatics tools for module extraction through network integration**

| Tool | URL | Refs |
|---|---|---|
| *Active-module detection through network projection of omics data* | | |
| jActiveModules | http://apps.cytoscape.org/apps/jactivemodules | 48 |
| MATISSE | http://acgt.cs.tau.ac.il/matisse | 165 |
| PinnacleZ | http://apps.cytoscape.org/apps/pinnaclez | 62 |
| GXNA | http://stat.stanford.edu/~serban/gxna | 52 |
| BioNet | http://bionet.bioapps.biozentrum.uni-wuerzburg.de | 166 |
| COSINE | http://cran.r-project.org/web/packages/COSINE/index.html | 104 |
| SANDY | http://sandy.topnet.gersteinlab.org | 81 |
| HotNet | http://ccmbweb.ccv.brown.edu/hotnet | 67 |
| PARADIGM | http://sbenz.github.com/Paradigm | 70 |
| MEMo | http://cbio.mskcc.org/memo | 73 |
| Multi-Dendrix | http://compbio.cs.brown.edu/software | 37 |
| RegMOD | http://www.biomedcentral.com/1471-2105/11/26/additional | 45 |
| NetWalk and FunWalk | http://netwalkersuite.org | 76 |
| ResponseNet | http://bioinfo.bgu.ac.il/respnet | 75 |
| ClustEx | http://www.mybiosoftware.com/pathway-analysis/5495 | 42 |
| SAMBA | http://acgt.cs.tau.ac.il/samba | 82 |
| cMonkey | http://bonneaulab.bio.nyu.edu/biclustering.html | 69 |
| COBRAv2.0 | http://opencobra.sourceforge.net/openCOBRA/Welcome.html | 85 |
| TieDIE | https://sysbiowiki.soe.ucsc.edu/tiedie | 167 |
| *Network comparisons across species to identify conserved modules* | | |
| PathBLAST | http://www.pathblast.org | 114 |
| NetworkBLAST | http://www.cs.tau.ac.il/~bnet/networkblast.htm | 168 |
| NetworkBLAST-M | http://www.cs.tau.ac.il/~bnet/License-nbm.htm | 116 |
| IsoRankN | http://groups.csail.mit.edu/cb/mna | 169 |
| Graemlin | http://graemlin.stanford.edu | 119 |
| NeXus | http://csbio.cs.umn.edu/neXus/help.html | 157 |
| Multi-species cMonkey | http://bonneaulab.bio.nyu.edu/biclustering.html | 158 |
| *Differential analysis of interaction networks to identify dynamic modules* | | |
| DDN | http://www.cbil.ece.vt.edu/software.htm | 170 |
| DNA | http://www.somnathdatta.org/Supp/DNA | 171 |
| *Integration of diverse types of interaction networks to identify composite modules* | | |
| PanGIA | http://prosecco.ucsd.edu/PanGIA | 147 |

Mitra *et al*, Nature Reviews Genetics (2013)

**Table 1 | Software tools based on network propagation**

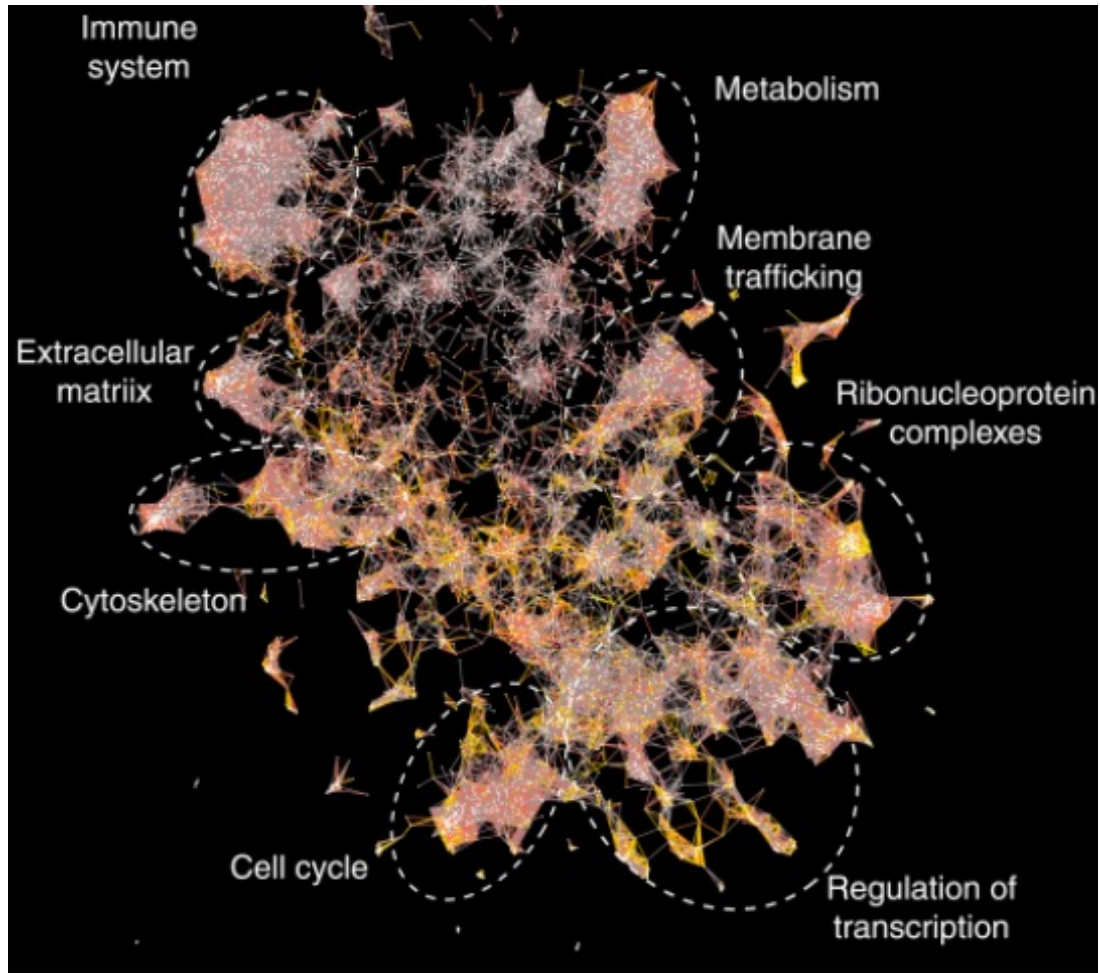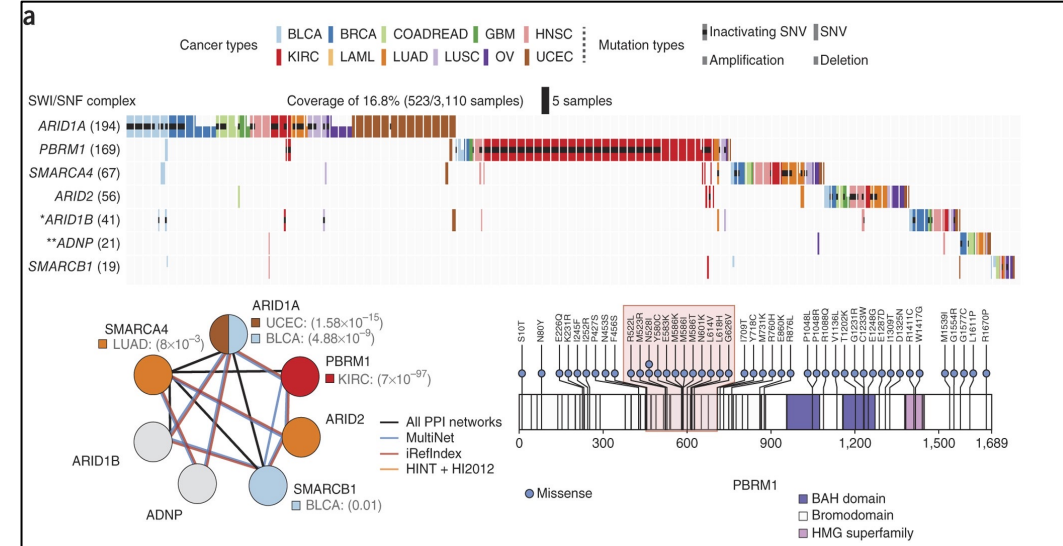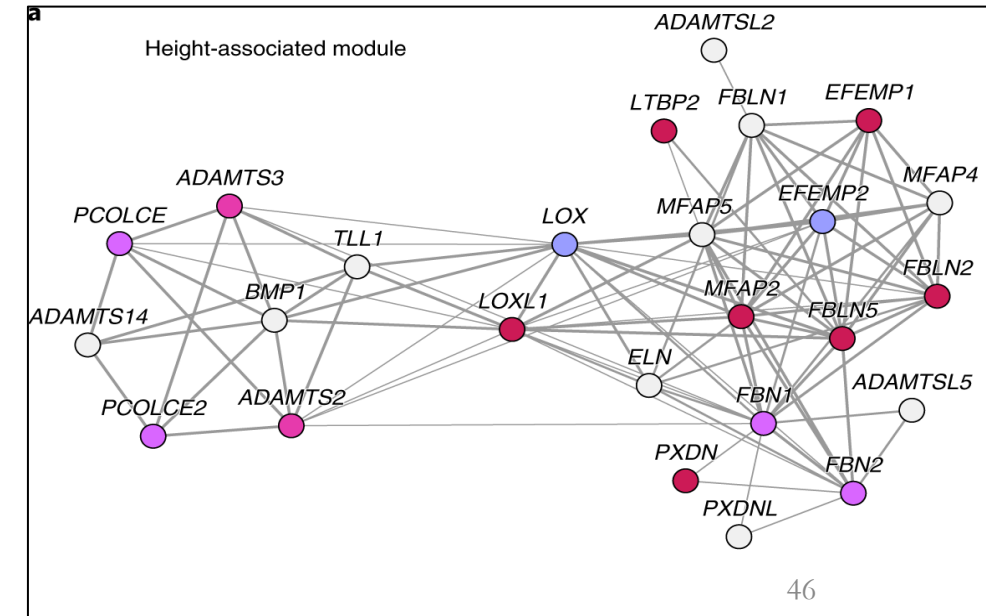| Tool | Goal | Type | Platform | Web site |
|---|---|---|---|---|
| *Function prediction* | | | | |
| DSD[48] and capDSD[34] | Function prediction | Single network | Web server and software for download | http://dsd.cs.tufts.edu/server/ and http://dsd.cs.tufts.edu/capdsd |
| GeneMANIA[103] | Function prediction | Single network | Cytoscape plugin | http://apps.cytoscape.org/apps/genemania |
| Mashup[56] | Function prediction | Integrative | Software for download | http://mashup.csail.mit.edu/ |
| RIDDLE[70] | Function prediction | Single network | Web server | http://www.functionalnet.org/RIDDLE/ |
| *Disease characterization* | | | | |
| CATAPULT[82] | Gene prioritization | Integrative | Web server and software for download | http://marcottelab/index.php/Catapult |
| Cytoscape 'diffuse' service[104] | General propagation | 1D and 2D | Software for download | • http://cytoscape.org • Native in version 3.5 and greater |
| DADA[80] | Gene prioritization | 1D | Software for download | http://compbio.case.edu/dada/ |
| Exome Walker[72] | Gene prioritization | 1D | Web server | http://compbio.charite.de/ExomeWalker |
| GUILD[105] | Gene prioritization | 1D | Software for download | http://sbi.imim.es/web/index.php/research/software/guildsoftware |
| HotNet2 (REF. 30) | Module detection | 2D | Software for download | http://compbio.cs.brown.edu/projects/hotnet2/ |
| NBS[89] | Patient stratification | Integrative | Software for download | http://chianti.ucsd.edu/~mhofree/NBS/ |
| NetQTL[79] | Gene prioritization and module detection | 1D | Software for download | https://www.ncbi.nlm.nih.gov/CBBresearch/Przytycka/index.cgi#netqtl |
| PRINCIPLE[106] | Gene prioritization and module detection | 1D | Cytoscape plugin | http://www.cs.tau.ac.il/~bnet/software/PrincePlugin/ |
| SNF[90] | Patient stratification | Integrative | Software for download | http://compbio.cs.toronto.edu/SNF/SNF/Software.html |
| TieDIE[91] | Module detection | Integrative | Software for download | https://sysbiowiki.soe.ucsc.edu/tiedie |
| ToppGene[107] | Gene prioritization | 1D | Web server | https://toppgene.cchmc.org/ |

Cowen *et al*, Nature Reviews Genetics (2017)

47

# Existing algorithms do not have rigorous, theoretical guarantees

Most algorithms assess their performance using real biological datasets, e.g.

- Runtime
- Recovering known biological findings
- Discovery of potentially new biological insights

But **do not** assess performance on generative model of the data

-> obscures fundamental issues shared across algorithms



High

Low

**Altered Subnetwork Problem:**
**Given**:
1) Network G = (V,E)
2) Vertex scores $X_v$ (usually derived from p-values)
**Goal**: Identify high-scoring subnetworks H of G

# Many algorithms output <u>**very large subnetworks**</u>



"Many algorithms are based on the score defined by jActiveModules [8], including PANOGA [9], dmGWAS [10], EW-dmGWAS [11], PINBPA [12], GXNA [13], and PinnacleZ [14]. Others, such as BioNet [15, 16] and Sig-Mod [17] are based on a score adapted to integer linear programming. These methods are also widely applied in the current literature [18, 19, 20, 21, 22, 14, 23, 24, 25, 26], even though the above approaches have been reported to consistently result in subnetworks that are large, and therefore difficult to interpret biologically [13, 27, 28]."

"Network module identification—a widespread theoretical bias and best practices" by Nikolayeva *et al* (Methods 2018)

**Altered Subnetwork Problem:**

**Given**:
1) Network G = (V,E)
2) Vertex scores $X_v$ (usually derived from p-values)

**Goal**: Identify high-scoring subnetworks H of G

**jActiveModules**/Cytoscape (Ideker et al, 2002): maximizes function over <u>connected subgraphs</u>

$$\arg\max_{\text{connected } S} \frac{1}{\sqrt{|S|}} \sum_{v \in S} X_v$$

49

# A simple simulation with an implanted subnetwork

Network has **10,000** vertices, implanted altered subnetwork A has **500** vertices
- Vertex scores in A are ~2 standard deviations larger than avg

jActiveModules outputs a subnetwork with **2505** vertices (5x increase!)



True

Estimated

High

Low

Many **heuristics** for reducing size – but effectiveness is unclear

# Our contributions

1. **Generative model** for altered subnetworks

2. Show issue of identifying large subnetworks is due to **statistical bias**

3. Develop NetMix algorithm which is **asymptotically unbiased**

Extensions:

- NetMix2 algorithm which uses network propagation (random walks)
- Anomaly detection in statistics/ML

High

Low

# Generative model: Altered Subnetwork Distribution

- G=(V, E) is a graph

- A ⊆ V is a <u>connected</u> subgraph, or the altered subnetwork



Vertex scores $(X_v)_{v \in V}$ are distributed as

$$X_v \sim \begin{cases} N(\mu, 1) & \text{if } v \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

**Altered Subnetwork Problem**: Given graph G and vertex scores $(X_v)_{v \in V}$
distributed as

$$X_v \sim \begin{cases} N(\mu, 1) & \text{if } v \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

find the altered subnetwork A.

**Altered Subnetwork Problem**: Given graph G and vertex scores $(X_v)_{v \in V}$ distributed as

$$X_v \sim \begin{cases} N(\mu, 1) & \text{if } v \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

find the altered subnetwork A.

**Hard** to solve ASP

**Easy** to solve ASP

**Impossible** to detect A



$\mu_{detect}$

Large $\mu$

Small $\mu$

**Altered Subnetwork Problem**: Given graph G and vertex scores $(X_v)_{v \in V}$ distributed as

$$X_v \sim \begin{cases} N(\mu, 1) & \text{if } v \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

find the altered subnetwork A.

**Impossible** to detect A

**Hard** to solve ASP

**Easy** to solve ASP

$\mu_{detect}$

Small $\mu$

Large $\mu$

**Altered Subnetwork Problem**: Given graph G and vertex scores $(X_v)_{v \in V}$ distributed as

$$X_v \sim \begin{cases} N(\mu, 1) & \text{if } v \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

find the altered subnetwork A.

Theorem: **Maximum Likelihood Estimator (MLE)** of the altered subnetwork A is:

$$\widehat{A}_{\text{MLE}} = \underset{\substack{S \subseteq V \\ S \text{ connected}}}{\text{argmax}} \left( \frac{1}{\sqrt{|S|}} \sum_{v \in S} X_v \right)$$

**Altered Subnetwork Problem**: Given graph G and vertex scores $(X_v)_{v \in V}$ distributed as

$$X_v \sim \begin{cases} N(\mu, 1) & \text{if } v \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

find the altered subnetwork A.

Theorem: **Maximum Likelihood Estimator (MLE)** of the altered subnetwork A is:

$$\widehat{A}_{\text{MLE}} = \underset{\substack{S \subseteq V \\ S \text{ connected}}}{\text{argmax}} \left( \frac{1}{\sqrt{|S|}} \sum_{v \in S} X_v \right)$$

**MLE = jActiveModules!**

jActiveModules paper (Ideker et al, 2002) does not describe generative model nor the connection to the MLE

# MLE is biased estimator

$$\text{Bias}\left(\frac{|\widehat{A}_{\text{MLE}}|}{n}\right) \triangleq E\left[\frac{|\widehat{A}_{\text{MLE}}|}{n}\right] - \frac{|A|}{n}$$

We observe that MLE has **positive bias:** MLE overestimates the size $\frac{|A|}{n}$ of the altered subnetwork on average (where n=|V|)



$|A|/n$

| | |
|---|---|
| —— | 0 |
| —— | 0.01 |
| —— | 0.05 ✹ (Simulation from before) |
| —— | 0.15 |
| —— | 0.25 |

# MLE is biased estimator

$$\text{Bias}\left(\frac{|\widehat{A}_{\text{MLE}}|}{n}\right) \triangleq E\left[\frac{|\widehat{A}_{\text{MLE}}|}{n}\right] - \frac{|A|}{n}$$

We observe that MLE has **positive bias:** MLE overestimates the size $\frac{|A|}{n}$ of the altered subnetwork on average (where n=|V|)



$|A|/n$

— 0

— 0.01

— 0.05 ✦ (Simulation from before)

— 0.15

— 0.25

**We prove** MLE is <u>asymptotically biased:</u> $\lim_{n \to \infty} \text{Bias}\left(\frac{|\widehat{A}_{\text{MLE}}|}{n}\right) > 0$

Assuming number of connected subgraphs is exponential
(RECOMB 2020; ICML 2021; unpublished w/ H Schmidt)

# How to reduce bias?

**Key idea:** Model the distribution of the vertex scores <u>before</u> using the network



Fit vertex scores to **Gaussian Mixture Model (GMM):**

$$X_v \sim (1 - \alpha) \cdot N(0, 1) + \alpha \cdot N(\mu, 1)$$

$\alpha$ = proportion of vertices in altered subnetwork
$\mu$ = mean of altered subnetwork distribution

# GMM yields less biased estimate of altered subnetwork size

**MLE:** $\widehat{A}_{\mathrm{MLE}} = \underset{\substack{S \subseteq V \\ S \text{ connected}}}{\mathrm{argmax}} \left( \dfrac{1}{\sqrt{|S|}} \sum_{v \in S} X_v \right)$

**vs**

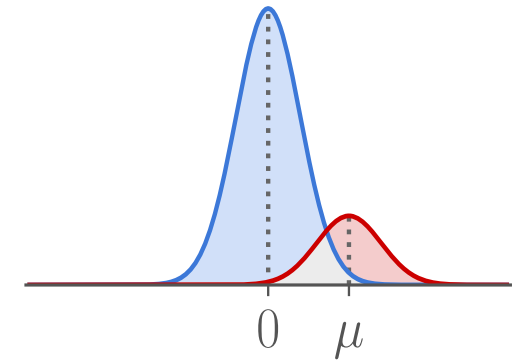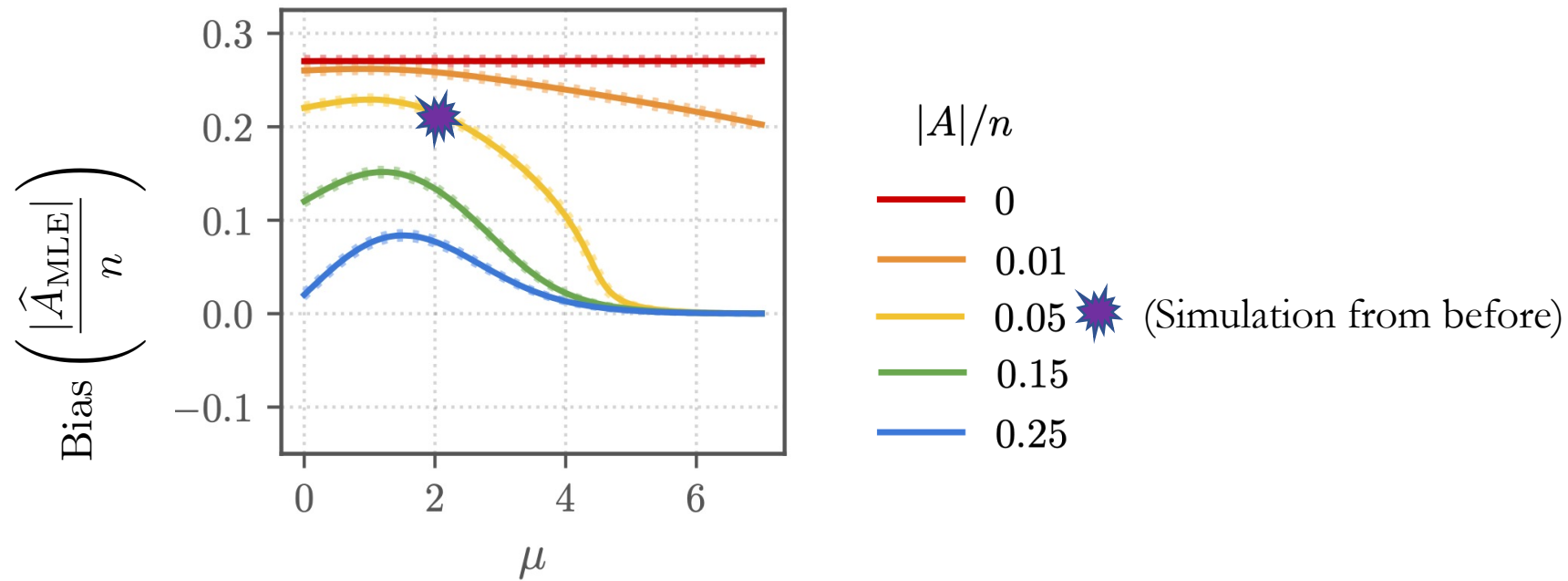**GMM:** Fit vertex scores $X_v$ to GMM

$$X_v \sim (1 - \alpha) \cdot N(0, 1) + \alpha \cdot N(\mu, 1)$$

and estimate GMM parameters $\quad \widehat{\alpha}_{\mathrm{GMM}}, \widehat{\mu}_{\mathrm{GMM}}$



**We prove:** GMM estimator is asymptotically unbiased (ICML 2021)

-> **Model mis-specification** helps!
(Fitting ASD with GMM)

$\alpha$ = proportion of vertices in altered subnetwork
$\mu$ = mean of altered subnetwork distribution

62

# NetMix Algorithm

Given vertex scores $(X_v)_{v \in V}$ and graph G:

1. Fit scores to GMM using EM, and compute *responsibilities* $r_v = P(v \in A \mid X_v)$
2. Find <u>connected subnetwork</u> $\widehat{A}_{\mathrm{NetMix}}$ with GMM-estimated size and largest total responsibility



Scores **X**

Low ▮▮▯▮▮ High

$N(0,1)$   $N(\widehat{\mu}_{\mathrm{GMM}}, 1)$

Fit scores to Gaussian
Mixture Model

$\widehat{A}_{\mathrm{NetMix}}$

Identify connected subnetwork with largest
total responsibility (ILP)

Interaction network
$G = (V, E)$

# Results – simulated data

Real PPI network (n ≈ 15,000 vertices)
Implanted altered subnetwork A
w/ size |A| = 0.05n = 750

NetMix — jActiveModules* — heinz (FDR = 0.001) — heinz (FDR = 0.1) — heinz (FDR = 0.5)

# Results – differential gene expression + somatic mutations in cancer



(A) (B) (C)

Legend: NetMix — jActiveModules* — heinz (FDR=.001) — heinz (FDR=.1) — heinz (Top |Â|) — Top |Â| $p$-values

157 gene expression experiments from Expression Atlas (Petryszak et al, 2015)

**Cancer driver gene** prediction:
- Using MutSigCV2 $p$-values and multiple interaction networks

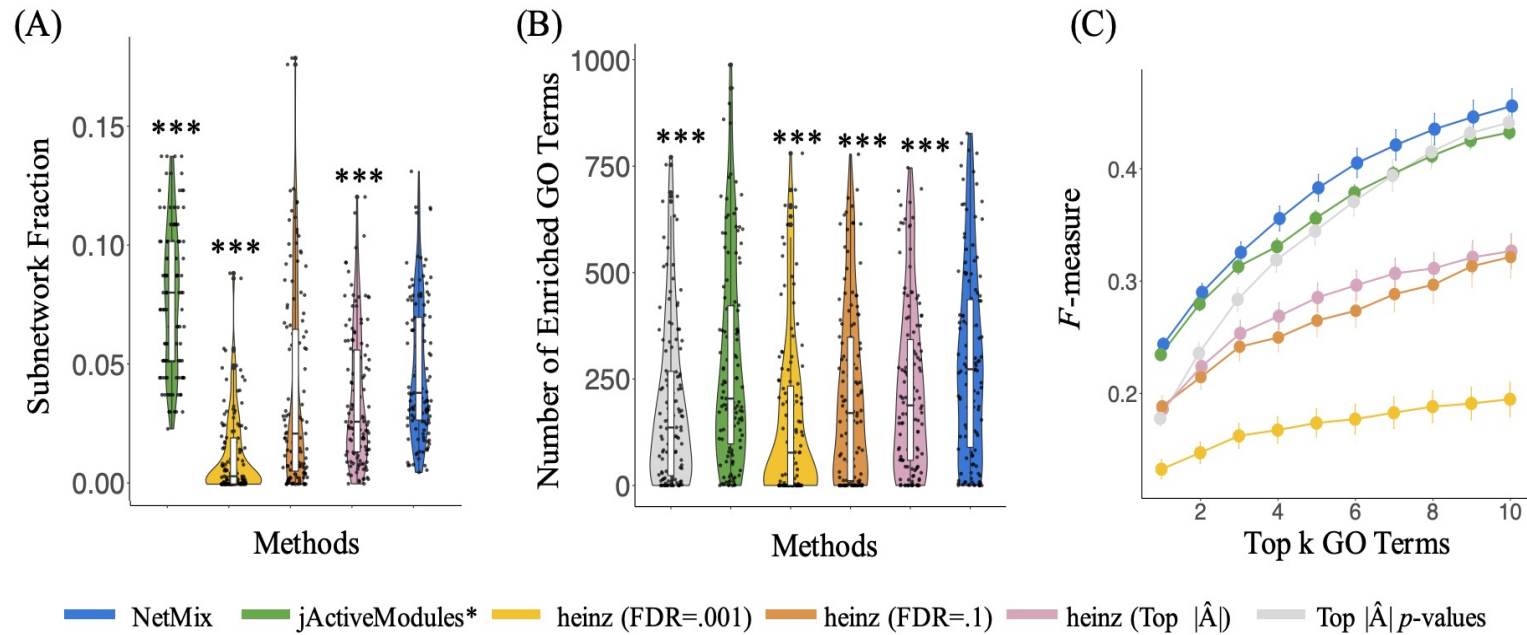| Method | Network | | | |
|---|---|---|---|---|
| | None | HINT+HI | iRefIndex | ReactomeFI |
| jActiveModules* | 2,136 / 0.155 | 1,575 / 0.191 | 1,815 / 0.174 | 557 / 0.261 |
| jActiveModules (Greedy search) | N.A. / N.A. | N.A. / N.A. | N.A. / N.A. | N.A. / N.A. |
| jActiveModules (Simulated annealing) | N.A. / N.A. | 12,284 / 0.086 | 15,046 / 0.074 | 8,329 / 0.118 |
| heinz (FDR = 0.001) | 115 / 0.205 | 119 / 0.216 | 109 / 0.217 | 114 / 0.215 |
| heinz (FDR = 0.1) | 259 / 0.244 | 249 / 0.264 | 259 / 0.255 | 253 / 0.215 |
| Hierarchical Hotnet | N.A. / N.A. | 228 / 0.214 | 297 / 0.215 | 228 / 0.214 |
| NetMix | 307 / **0.254** | 263 / **0.277** | 296 / **0.270** | 264 / **0.270** |

# NetMix2: extension to other distributions and graph topologies

# Anomaly detection in statistics/ML

Normal means problem: Data $X_1, \ldots, X_n$ independently distributed as

$$X_i \sim \begin{cases} N(\mu, 1) & \text{if } i \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

where anomaly $A \in \mathcal{S}$ is a member of anomaly family $\mathcal{S}$

Anomalies are <u>connected subgraphs</u>

$\boxed{\mathcal{S} = \mathcal{C}_G}$

$A$

Connected Subgraphs

Anomalies are <u>time intervals</u>

$\boxed{\mathcal{S} = \mathcal{I}_n}$

$A$

Intervals

Anomalies are <u>submatrices</u>

$\boxed{\mathcal{S} = \mathcal{M}_N}$

columns of $A$

rows of $A$

Submatrices

ICML 2021: **We extend theoretical results** and show: MLE is biased iff number of sets in anomaly family $\mathcal{S}$ containing A is <u>exponential</u>

# MLE is biased for spatial anomalies

**Spatial adjacency graph**
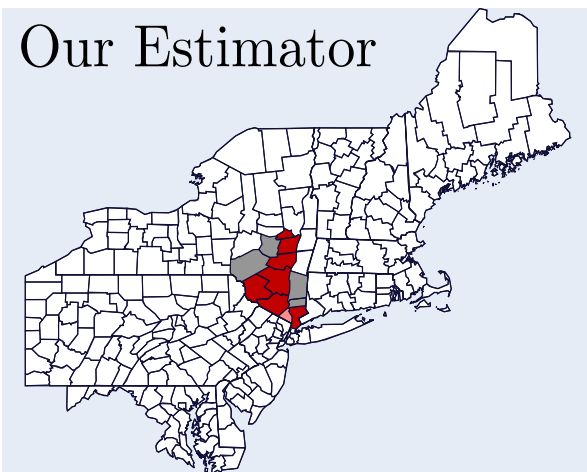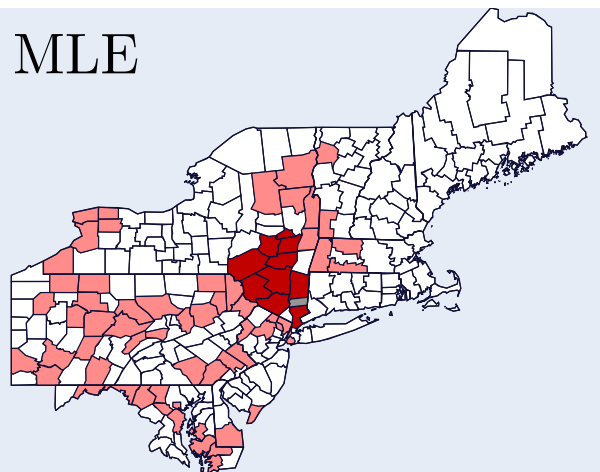


High score

Low score

- Vertices = points in space
- Edges connect adjacent points in space
- Score = disease incidence



A **spatial scan statistic**

M Kulldorff - Communications in **Statistics**-Theory and methods, 1997 - Taylor & Francis

The **scan statistic** is commonly used to test if a one dimensional point process is purely random, or if any clusters can be detected. Here it is simultaneously extended in three directions:(i…

☆ Save  🗩 Cite   Cited by 4448   Related articles   All 8 versions   Web of Science: 2406

MLE = network version of widely-used *"spatial scan statistic"*

MLE                        Our Estimator



MLE                    Our Estimator



■ True positive    ■ False positive    ■ False negative

**Simulated disease outbreak in northeast US**

**Real data: NYC breast cancer incidence**

# Summary

**Generative model** for altered subnetworks

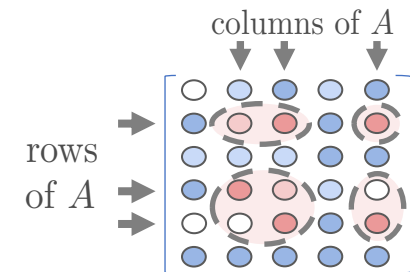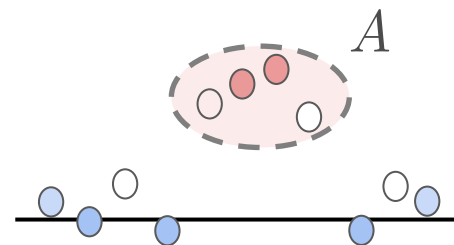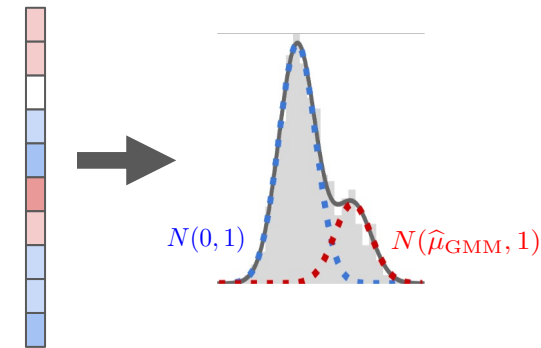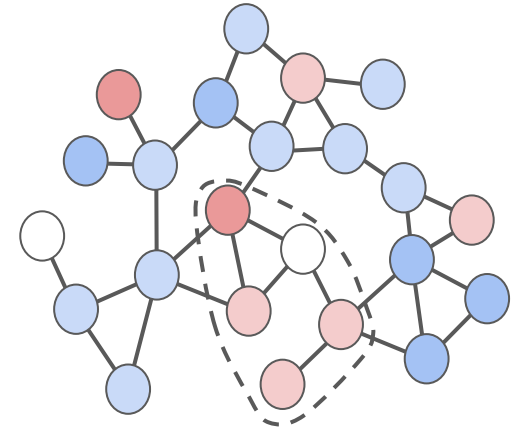**We show:** MLE is asymptotically biased for connected subgraphs

**NetMix/NetMix2**: asymptotically unbiased altered subnetwork algorithms

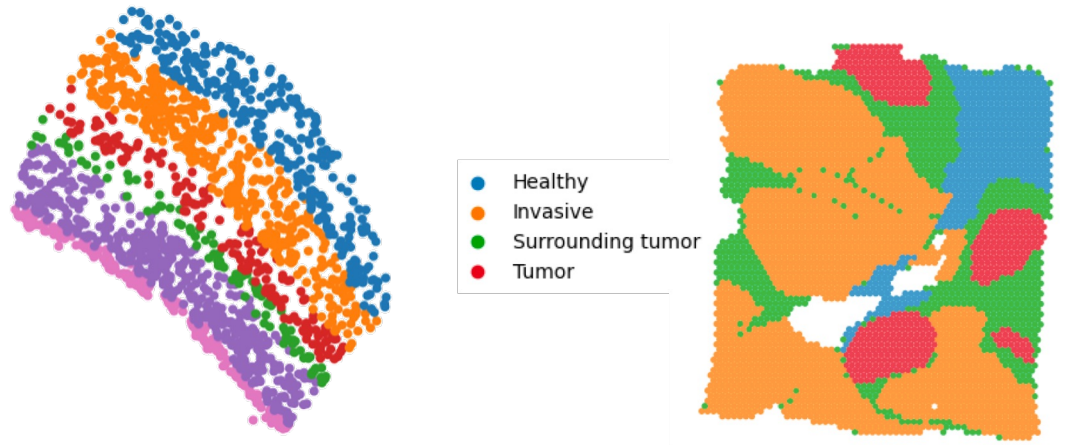Idea: fit vertex scores to mixture model **before** using network

Results extend to **anomaly detection** in machine learning/statistics

Future idea: anomaly detection in spatial transcriptomics?



$N(0,1)$ $N(\hat{\mu}_{\mathrm{GMM}}, 1)$

$A$

columns of $A$

rows of $A$

# My thesis: computational methods for understanding complex biological systems

## Spatial biology



Healthy
Invasive
Surrounding tumor
Tumor

<u>Spatial variation in gene expression</u>
- Ma*, **Chitra***, et al. *RECOMB 2022 + Cell Systems.*
- **Chitra** et al. *RECOMB 2024 + in review at Nature Methods.*

<u>Cell-cell interactions</u>
- Sarkar*, **Chitra***, et al. *In submission at ISMB 2024.*

## Network interactions and anomalies



<u>Altered subnetwork identification</u>
- Reyna*, **Chitra***, et al. *RECOMB 2020 + JCB.*
- **Chitra** et al. *ICML 2021.*
- **Chitra***, Park*, Raphael. *RECOMB 2022 + JCB.*

<u>Learning genetic interactions</u>
- **Chitra***, Arnold*, Raphael. *In review at Nature Genetics.*
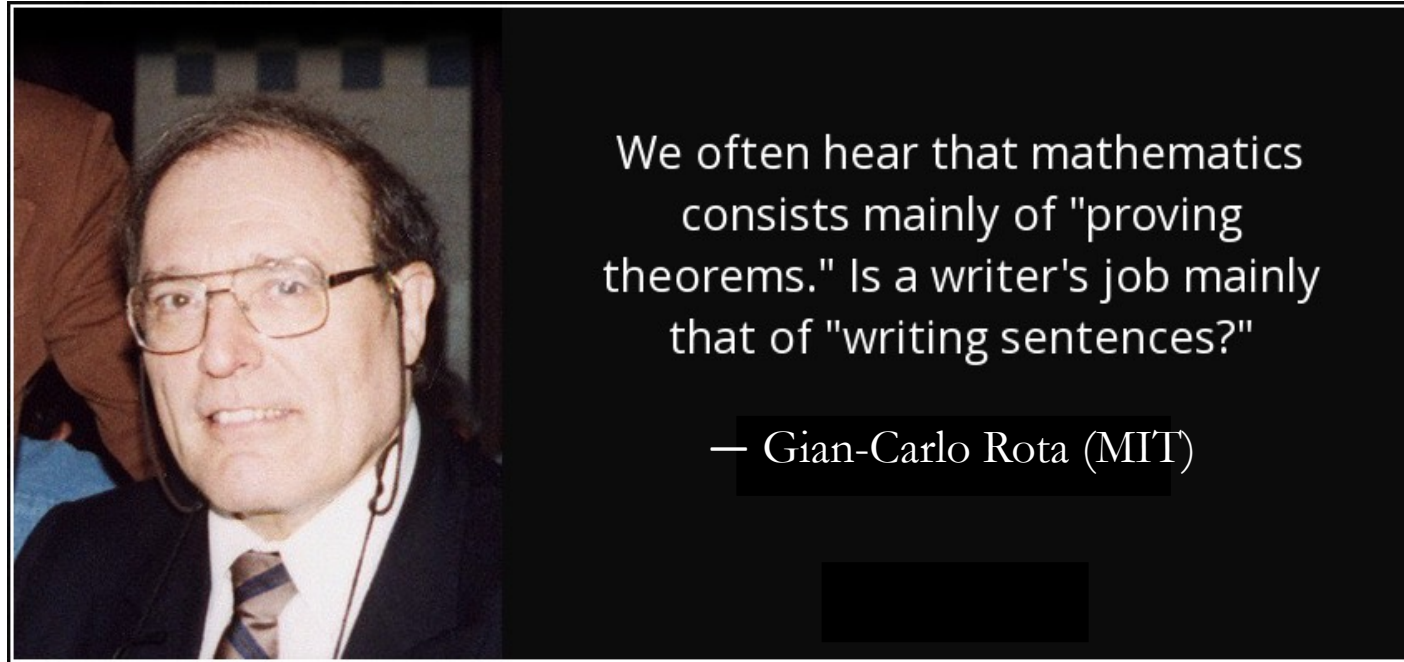- Shuaibi*, **Chitra***, Raphael. *In submission at RECOMB-CCB.*

<u>Machine learning + data mining</u>
- **Chitra** and Raphael. *ICML 2019.*
- **Chitra** and Musco. *WSDM 2020.*

70

* indicates joint first authorship

# What does a computational biology researcher do?



We often hear that mathematics consists mainly of "proving theorems." Is a writer's job mainly that of "writing sentences?"

— Gian-Carlo Rota (MIT)

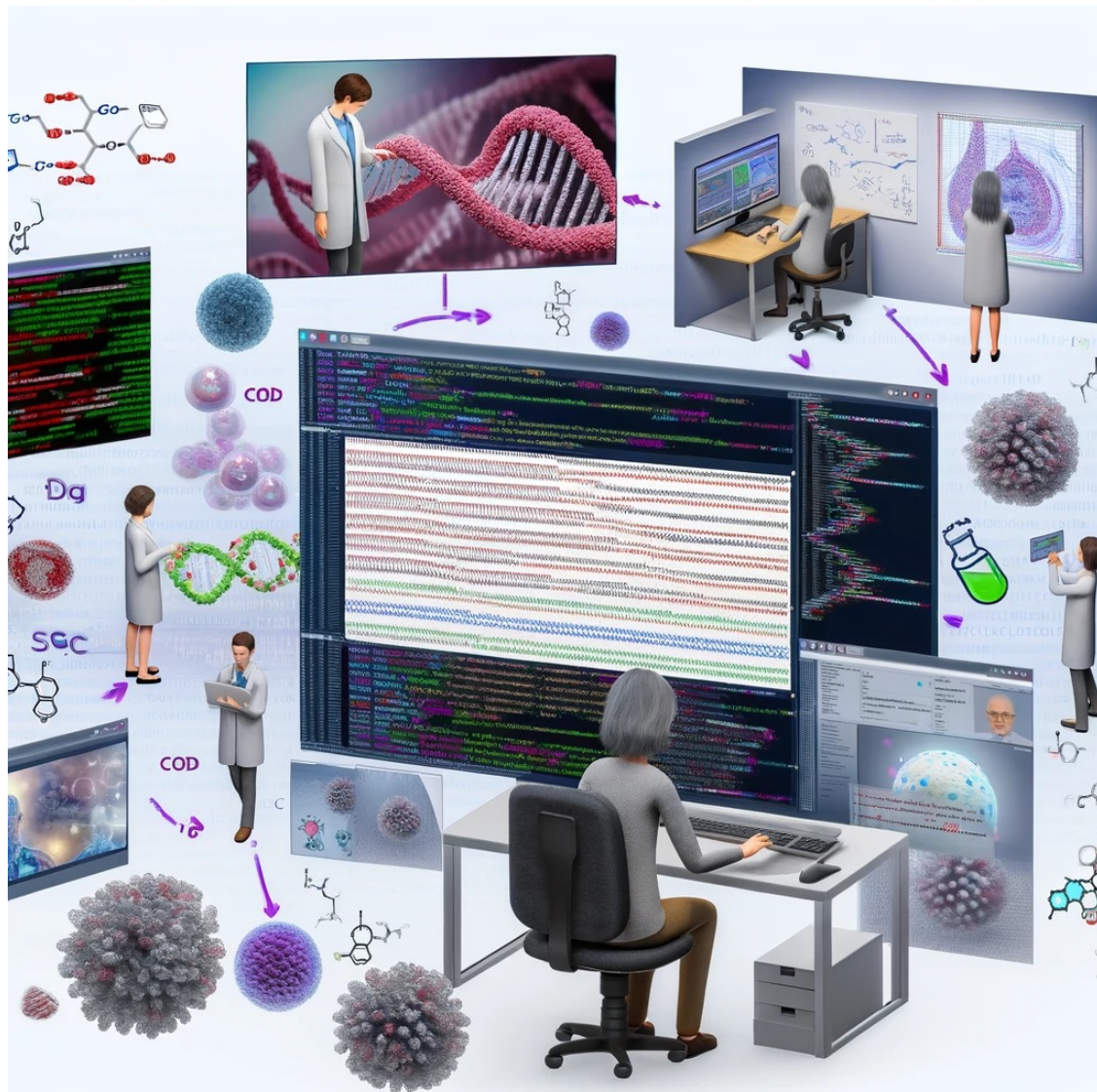## Is a computational biologist's job mainly that of "analyzing biological data?"

No! We also try to (1) identify biologically interesting problems and (2) find mathematically "elegant" solutions to problems

UT **You**

Generate an image answering "what does a computational biologist do?"

# Thank you!



**Advisor**

Ben Raphael

**Committee**

Bernard Chazelle

Ellen Zhong

Yuri Pritykin

Fei Chen

# Almost a decade working with Ben



Group retreat, summer 2015



Group retreat, summer 2023

# Acknowledgments

**Collaborators/co-authors:**

Ben Raphael
Matt Reyna
Rebecca Elyanow
Kimberly Ding
Jasper Lee
Tyler Park
Cong Ma
Shirley Zhang
Brian Arnold
Sereno Lopez-Darwin
Hirak Sarkar
Kohei Sanno
Ahmed Shuaibi
Julian Gold
Clover Zheng
Sunay Joshi
Chris Musco
Tarun Chitra

**Raphael group (past and present):**

Simone Zaccaria
Ron Zeira
Pijus Simonaitis
Gryte Satas
Matt Myers
Hongyu Zheng
Palash Sashittal
Uyen Mai
Metin Balaban
Richard Zhang
Alexander Strzalkowski
Henri Schmidt
Xinhao Liu
Akhil Jakatdar
Gary Hu
Peter Halmos
Gillian Chu
Clover Zheng
Maya Gupta
Madelyne Xiao

+ support of numerous friends + family



GRFP
NSF GRADUATE RESEARCH FELLOWSHIP PROGRAM

NIH National Human Genome Research Institute

Siebel Scholars

# Extra content

# Belayer – simulated data



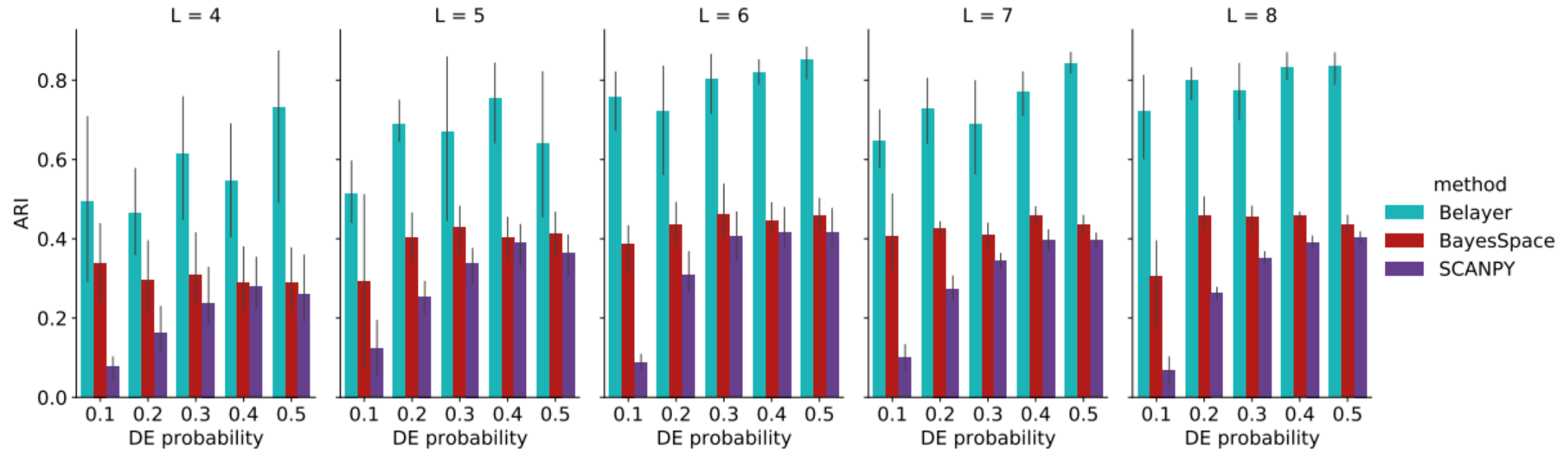Figure S6: Comparison of Belayer, BayesSpace, and SCANPY in identifying spatially distinct cell clusters in the second simulation. Performance of each method is evaluated according to the Adjusted Rand Index (ARI) and shown for different values of the number $L$ of layers and differential expression (DE) probability. Error bars indicate variation from 5 randomly simulated datasets for each parameter setting.

# Belayer accurately identifies cortical layers in human DLPFC



Takeaway: **Global** spatial model (Belayer) > **local** models

# Different accuracy metrics: DLPFC

Figure S9: Comparison of Belayer and other methods from Figure 3 with different metrics.

# Belayer identifies spatially coherent cortical layers

(mouse somatosensory cortex data, SlideSeqV2)



SpaGCN: reference-free + **local** spatial model (GNN)

Belayer: reference-free + **global** spatial model

Approximate layer boundaries

Layer depth

RCTD cell types

SpaGCN

Belayer

RCTD cell type labels

- L2/3 IT
- L4
- L5 IT
- L5 PT
- L6 CT
- L6 IT
- L6b

ARI=0.34

ARI=0.37

Data from Stickels et al., 2021
RCTD [Cable et al., 2021]

# Belayer identifies spatially varying genes
(mouse somatosensory cortex data, SlideSeqV2)

**A** Spatially resolved transcriptomics data

Coordinates

Gene expression

**GASTON**
Gradient Analysis of Spatial Transcriptomics Organization with Neural networks

**Topographic map** of tissue slice

Isodepth

**B**

Input: Spatial coordinates

$x$ $y$

Hidden layer: **isodepth**

$d$

$a_1$ $a_2$ $a_3$ $a_4$ $\cdots$ $a_G$

Output: Gene expression

Downstream analyses
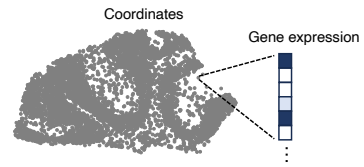
**C** Spatial domains

- Domain 1
- Domain 2
- Domain 3

**D** Continuous gradients and discontinuous variation in gene expression

Expression

Isodepth

**E** Spatial variation in cell type organization

Cell type proportion

Cell type 1
Cell type 2
Cell type 3
Cell type 4
Cell type 5

Isodepth

**F** Continuous gradients in tumor microenvironment

Cancer

Stroma

Expression

Isodepth

82

# GASTON – model selection (elbow)

# GASTON – cerebellum (spatial domains)

GASTON – cerebellum (spatial expression patterns)

# GASTON – colorectal tumor

**A** **B** **C**

Domain 1 (tumor)
Domain 2 (tumor-adjacent stroma)
Domain 3
Domain 4
Domain 5

**D** Intrastromal variation / Discontinuity / Intratumoral variation

Stroma / Tumor

564 | 56 | 93
50 | 67
186 | 556

**E**

OXIDATIVE PHOSPHORYLATION
CHOLESTEROL HOMEOSTASIS
EPITHELIAL MESENCHYMAL TRANSITION
HYPOXIA
MYOGENESIS
COAGULATION
COMPLEMENT
ESTROGEN RESPONSE EARLY
ESTROGEN RESPONSE LATE
P53 PATHWAY
TNFA SIGNALING VIA NFKB

GeneRatio 0.2 0.3
p.adjust 0.006 0.004 0.002

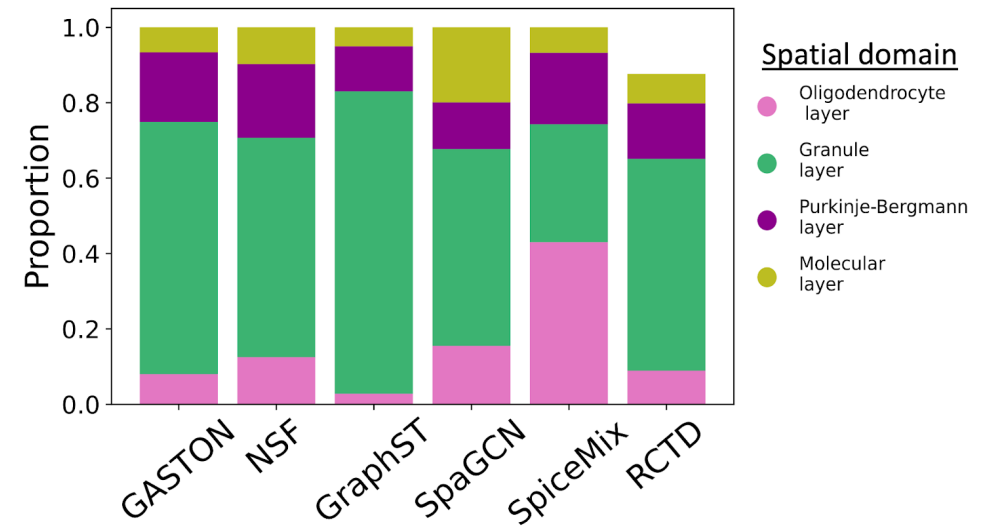Intrastromal variation
Discontinuity
Intratumoral variation

Type I II III

**F** $COX7B$

**G** $SCD$

**H** $ACTA2$

**I** $TAGL$

**J** $COL1A2$ expression

**K** $COL1A2$

**L** GASTON $COL1A2$ expression function

## Seurat cell types

CellType
Fibroblast
Tumor
Lamina Propria

## Tumor growth rate

**A** *TAGLN* expression

**B** *ACTA2* expression

GASTON – colorectal tumor (more patterns)

**D** *LGR5* expression

**C** *LGR5*

**E** *THBS1*

**F** *FUCA1*

# Comparison of domains on colorectal tumor

# Olfactory bulb (Stereo-seq)

9,825 spots × 27,106 genes

- ● Olfactory nerve layer (ONL)
- ● Glomerular layer (GL)
- ● External plexiform layer (EPL)
- ● Mitral cell layer (MCL)
- ● Internal plexiform layer (IPL)
- ● Granule cell layer (GCL)
- ● Rostral migratory stream (RMS)

**DAPI Stain**



**Isodepth and spatial gradients**



**SpaceFlow** (diffusion pseudotime)



**GASTON**



**SpaGCN**



**Cell type-*attributable* gradient**

**Other attributable gradient**



Cell type proportion vs Isodepth

- Astro
- Granule
- Mes
- Mitral/Tufted
- OEC
- OSN
- Periglomerular
- Transition

$GAD1_{, all}$ *cell types* — log(pooled CPM) vs Isodepth

*Granules only*

*Astrocytes only*

89

# Comparison b/w GASTON and Belayer



Cerebellum (Slide-seqV2) | Colorectal cancer tumor (10x Visium) | Olfactory bulb (Stereo-seq)

GASTON

Belayer (unsupervised)

**GASTON – DLPFC**

**A**

ARI

GASTON  Belayer  SpaGCN  STAGATE

**B**

Manual annotation    GASTON    Belayer

- ■ DLPFC/GASTON/Belayer Layer 1
- ■ DLPFC Layer 2
- ■ DLPFC/GASTON/Belayer Layer 3
- ■ DLPFC Layer 4
- ■ DLPFC/GASTON/Belayer Layer 5
- ■ DLPFC/GASTON/Belayer Layer 6
- ■ DLPFC/GASTON/Belayer White Matter (WM)

**C**

Isodepth and spatial gradients

**D**

Correlation with isodepth

Belayer relative depth    GLM-PC1    PC1

# GASTON – SpaceFlow comparison (cerebellum)

**A**

**GASTON**



**B**

**SpaceFlow**
(diffusion pseudotime)



**C**

GASTON
Granule layer



**D**

SpaceFlow
Granule layer



**E**



Quartile coefficient of dispersion

- GASTON
- SpaceFlow

Oligodendrocyte layer, Granule layer, Purkinje-Bergmann layer, Molecular layer

# GASTON – mouse primary motor cortex (MERFISH)



**A** GASTON topographic map

**B** GASTON spatial domains

**C** ENVI "pseudo-depth" coordinate

L2/3

L5/6

L5

L6

**D** Acta2

**E** Chn2

**F** *Chn2* expression

**G** *Chn2* GASTON expression function

# GASTON – breast cancer (10x Xenium)

A — Chen et al (Spateo), GASTON, STAGATE, SpaGCN, GraphST

Legend:
- AGM
- Brain
- Branchial arch
- Cavity
- Connective tissue
- Dermomyotome
- Heart
- Liver
- Mesenchyme
- Neural crest
- Notochord
- Sclerotome

B

C — GASTON domain 7 (heart)

D — Top GO terms for 128 genes with continuous variation in GASTON domain 7 (heart)

- Heart Contraction (GO:0060047)
- Cardiac Muscle Contraction (GO:0060048)
- Striated Muscle Contraction (GO:0006941)
- Myofibril Assembly (GO:0030239)
- Cardiac Muscle Tissue Morphogenesis (GO:0055008)
- Sarcomere Organization (GO:0045214)
- Cardiac Ventricle Morphogenesis (GO:0003208)
- Heart Development (GO:0007507)
- Actomyosin Structure Organization (GO:0031032)
- Ventricular Cardiac Muscle Tissue Morphogenesis (GO:0055010)
- Ventricular Cardiac Muscle Tissue Development (GO:0003229)
- Regulation Of Heart Contraction (GO:0008016)
- Muscle Contraction (GO:0006936)
- Response To Muscle Stretch (GO:0035994)
- Cardiac Muscle Cell Development (GO:0055013)

-log10(Adjusted P-value)

E — Bmp4

F — Cacna1c

GASTON – mouse embryo day 9.5 (Stereo-seq)

95

# Application of GASTON to metabolomics (Clover Zheng)



Testing GASTON w/ simulated hexagonal geometries:

True isodepth d(x,y)

GASTON-estimated isodepth

# GMM yields less biased estimate of altered subnetwork size

**MLE:** $\widehat{A}_{\mathrm{MLE}} = \underset{\substack{S \subseteq V \\ S \text{ connected}}}{\mathrm{argmax}} \left( \frac{1}{\sqrt{|S|}} \sum_{v \in S} X_v \right)$

**vs**

**GMM:** Fit vertex scores $X_v$ to GMM

$$X_v \sim (1 - \alpha) \cdot N(0, 1) + \alpha \cdot N(\mu, 1)$$

and estimate GMM parameters $\widehat{\alpha}_{\mathrm{GMM}}, \widehat{\mu}_{\mathrm{GMM}}$

**We prove** (ICML 2021): GMM yields asymptotically unbiased estimates of $\alpha$, $\mu$, i.e.

$$\lim_{n \to \infty} |\widehat{\alpha}_{\mathrm{GMM}} - \alpha| = 0$$

$$\lim_{n \to \infty} |\widehat{\mu}_{\mathrm{GMM}} - \mu| = 0$$

Model mis-specification helps!
(Fitting ASD with GMM)

$\alpha$ = proportion of vertices in altered subnetwork

$\mu$ = mean of altered subnetwork distribution

# Challenge: Connectivity is a weak topological constraint!

Networks have small diameter – most subnetworks are "almost connected"

Algorithms not much better compared to not using interaction network



Simulations from our generative model where altered subnetwork is **connected subgraph**

# Network propagation (network diffusion)

Use of <u>random walks</u> to "propagate"/smooth vertex scores across network



High

Low

Network propagation

a

Direct neighbour    Shortest path    Network propagation

## Network propagation: a universal amplifier of genetic associations

Lenore Cowen, Trey Ideker, Benjamin J. Raphael & Roded Sharan ✉

_Nature Reviews Genetics_ **18**, 551–562 (2017) │ Cite this article

**18k** Accesses │ **257** Citations │ **41** Altmetric │ Metrics

# Network propagation uses <u>global</u> network structure



**a**   t=0        t=1        t=2        t=3    ...   t=∞

Cowen et al (Nature Reviews Genetics 2017)

Network propagation = Matrix-vector multiplication



Random walk
similarity matrix

Vertex scores

| Name | Similarity matrix |
|------|-------------------|
| Random walk | $W^k$ |
| Random walk with restart | $\alpha(I-(1-\alpha)W)^{-1}$ |
| Diffusion kernel | $e^{-\alpha W}$ |

Cowen et al (Nature Reviews Genetics 2017)

# Network propagation is standard for <u>ranking vertices</u>



Known

Unknown

High

Low

Rank vertices based on similarity to vertices w/ <u>known</u> characteristics e.g. genes associated with a specific disease (<u>binary</u> vertex scores $X_v$)



Random walk similarity matrix

Vertex scores

Personalized PageRank is **asymptotically optimal** for ranking in random graph models (PNAS 2017)

# How to use <u>network propagation</u> to identify altered subnetworks?



Network propagation

High

Low

**Question**: how to identify altered subnetwork from propagated gene scores?

# Existing network propagation methods use <u>ad hoc heuristics</u> to identify altered subnetworks



High

Low

Network propagation

**HotNet2**

**Pan-cancer network analysis identifies combinations of rare somatic mutations across pathways and protein complexes**

Mark D M Leiserson, Fabio Vandin, Hsin-Ta Wu, Jason R Dobson, Jonathan V Eldridge, Jacob L Thomas, Alexandra Papoutsaki, Younhun Kim, Beifang Niu, Michael McLellan, Michael S Lawrence, Abel Gonzalez-Perez, David Tamborero, Yuwei Cheng, Gregory A Ryslik, Nuria Lopez-Bigas, Gad Getz, Li Ding & Benjamin J Raphael ✉

*Nature Genetics* **47**, 106–114 (2015) | Cite this article

**39k** Accesses | **500** Citations | **122** Altmetric | Metrics

**PRINCE**

**Associating Genes and Protein Complexes with Disease via Network Propagation**

Oron Vanunu co, Oded Magger co, Eytan Ruppin, Tomer Shlomi, Roded Sharan ✉

Published: January 15, 2010 • https://doi.org/10.1371/journal.pcbi.1000641

Ex: **PRINCE**: "We aim at inferring <u>densely connected protein complexes that contain</u> <u>high scoring proteins</u> … we start with the top 100 [propagated] scoring proteins as seeds … <u>To each seed we iteratively add a neighboring protein</u> with the highest score … A <u>refinement phase</u> takes place where <u>proteins are removed</u> from a putative complex to ensure that … its member proteins are densely interacting."

**Issue:** These algorithms lack **<u>rigorous statistical guarantees</u>** – hard to investigate fundamental issues like bias

# Recent work shows existing approaches biased towards "high centrality" vertices

Algorithms benchmark against existing network algorithms – can hide <u>biases shared across methods</u>

**DOMINO: a network-based active module identification algorithm with reduced rate of false calls**

Hagai Levi, Ran Elkon (iD), Ron Shamir (iD) ✉

<u>Author Information</u>

Molecular Systems Biology (2021) 17: e9593 | https://doi.org/10.15252/msb.20209593

*"Our study reports on a different bias that is prevalent in AMI solutions: **<u>their tendency to report non-specific GO terms</u>**. …we observed that many enriched GO terms also appear on permuted datasets, suggesting that such enrichment stems from some proprieties of the network, algorithm, or the data that bias the results."*

**On the limits of active module identification**

Olga Lazareva, Jan Baumbach, Markus List, David B Blumenthal ✉    Author Notes

*Briefings in Bioinformatics*, Volume 22, Issue 5, September 2021, bbab066,

https://doi.org/10.1093/bib/bbab066

Published: 29 March 2021    Article history ▾

*"Our results indicate that classical but also supposedly bias-aware [altered subnetwork algorithms] extract disease modules **<u>based on the node degree</u>**"*

# Our work:

- Extend altered subnetwork generative model
  - Model different altered subnetwork topologies ("**subnetwork families**")
  - Derive <u>propagation family</u> – "approximates" subnetworks found by network propagation

- **<u>NetMix2</u>** algorithm for altered subnetwork identification with different subnetwork families
  - w/ propagation family: principled network propagation algorithm for altered subnetwork identification

- Simple baselines for evaluating network algorithms – *"scores only"* and *"network only"*

# Generative model: Altered Subnetwork Distribution



- G=(V, E) is interaction network
- $\mathcal{S}$ is **subnetwork family** (set of subsets of V)
- $A \in \mathcal{S}$ is the altered subnetwork

Vertex scores $(X_v)_{v \in V}$ are distributed as

$$X_v \sim \begin{cases} \mathcal{D}_a, & \text{if } v \in A, \\ \mathcal{D}_b, & \text{otherwise} \end{cases}$$

$\mathcal{D}_a$ = altered distribution (unknown)

$\mathcal{D}_b$ = background distribution (typically known)

---

**Example of distributions**: z-scores

$$\mathcal{D}_a = N(\mu, 1)$$
$$\mathcal{D}_b = N(0, 1)$$



---

**Examples of subnetwork families:**

Connected family   $\mathcal{S} = \mathcal{C}_G$ = connected subgraphs S ⊆ V

Edge-dense family   $\mathcal{S} = \mathcal{E}_{G,p}$ = subgraphs with density(S) > p

Cut family   $\mathcal{S} = \mathcal{T}_{G,\rho}$ = subgraphs with cut(S) < ϱ

# Generative model: Altered Subnetwork Distribution

- G=(V, E) is interaction network
- $\mathcal{S}$ is **subnetwork family** (set of subsets of V)
- $A \in \mathcal{S}$ is the altered subnetwork

Vertex scores $(X_v)_{v \in V}$ are distributed as

$$X_v \sim \begin{cases} \mathcal{D}_a, & \text{if } v \in A, \\ \mathcal{D}_b, & \text{otherwise} \end{cases}$$

High

Low

$A$

**Altered Subnetwork Problem (ASP)**: Given graph G, subnetwork family $\mathcal{S}$ and vertex scores $(X_v)_{v \in V}$, find altered subnetwork A.

ASP = <u>estimating parameters of distribution</u>

$|\mathcal{S}|$ is *large*

**Hard** to solve ASP

"Sweet spot" – both anomaly family and data are important

**Easy** to solve ASP without anomaly family

**Hard** to solve ASP

*Small* distance between distributions $\mathcal{D}_a, \mathcal{D}_b$

*Large* distance between distributions $\mathcal{D}_a, \mathcal{D}_b$

**Easy** to solve ASP without much info from data

$|\mathcal{S}|$ is *small*

Data $X_1, \ldots, X_n$ distributed as

$$X_i \sim \begin{cases} \mathcal{D}_a, & \text{if } i \in A, \\ \mathcal{D}_b, & \text{otherwise} \end{cases}$$

where anomaly $A \in \mathcal{S}$ is a member of anomaly family $\mathcal{S}$

109

**Hard** to solve ASP

$|\mathcal{S}|$ is *large*

"Sweet spot" - both network (subnetwork family) and vertex scores are necessary

**Hard** to solve ASP

**Easy** to solve ASP without network

*Small* distance between distributions $\mathcal{D}_a, \mathcal{D}_b$

*Large* distance between distributions $\mathcal{D}_a, \mathcal{D}_b$

**Easy** to solve ASP without vertex scores

**Important** to compare against naïve baselines that use
(1) "scores only"
(2) "network only"

$|\mathcal{S}|$ is *small*

110

# Propagation family

$$\mathcal{S} = \mathcal{M}_{\delta,p} : \text{Subgraphs S with} \quad \underline{M_{u,v} \geq \delta} \quad \text{for p fraction of } (u, v) \in S$$

Vertices are "close"
via random walk

(also require $M_{v,u} \geq \delta$ if M is not symmetric, eg personalized PageRank)

In RECOMB 2022 paper: some theory and simulations show <u>propagation family approximates subnetworks found by network propagation methods</u>

**Alternatively**: edge-dense subnetworks of *"similarity threshold graph"*



Vertex scores $X_v$    **Random walk similarity matrix M**    Vertex scores $X_v$

network propagation

Interaction network G

Similarity threshold graph $G_\delta$

# Simulations: Propagation family corresponds to the subnetworks identified by network propagation



Scores only = {vertices w/ top-|A| scores}

Network only = {vertices w/ top-|A| vertex centrality}

Network propagation = {vertices w/ top-|A| propagated scores}

$G$ = HINT+HI interaction network with $|G| \approx 15000$ nodes (Leiserson et al 2015)

Altered subnetwork A of size $|A|$=0.01n selected uniformly at random from subnetwork family $\mathcal{S}$

# Results: somatic mutations in cancer

NetMix2 outperforms other methods at identifying previously reported driver mutations in cancer.

| Method | Subnetwork size | STRING network | | | | | |
| | | CGC | | OncoKB | | TCGA | |
| | | Number | F-measure | Number | F-measure | Number | F-measure |
| --- | --- | --- | --- | --- | --- | --- | --- |
| NetMix2 | 280 | 132 | **0.3** | 133 | **0.313** | 151 | **0.546** |
| NetMix | 313* | 129 | 0.282 | 130 | 0.295 | 147 | 0.502 |
| Heinz (FDR=0.01) | 335 | 139 | 0.297 | 138 | 0.306 | 156 | 0.513 |
| NetSig | 773 | 145 | 0.211 | 172 | 0.257 | 84 | 0.161 |
| Hierarchical HotNet | 246 | 73 | 0.172 | 70 | 0.172 | 74 | 0.285 |
| Network Propagation | 280 | 86 | 0.195 | 89 | 0.210 | 98 | 0.354 |
| Scores-only | 280 | 126 | 0.286 | 127 | 0.3 | 145 | 0.524 |
| Network-only | 280 | 77 | 0.175 | 83 | 0.196 | 55 | 0.199 |

G = STRING protein interaction network

Vertex scores $X_v$ = MutSIg2CV z-scores computed based on frequency of somatic mutations in TCGA tumor samples

**Note:** "Scores-only" has good performance – how helpful is interaction network?

# Results: GWAS

Recent study by Carlin et al (iScience 2019) – evaluates how well methods identify <u>known disease reference genes</u>



Network propagation

"scores only"

# Results: GWAS

Recent study by Carlin et al (iScience 2019) – evaluates how well methods identify known <u>disease reference genes</u>



**Issue:** AUROC is poor metric for small reference sets! (<1% of 15,000 genes)

Network propagation

"scores only"

AUPRC

**Schizophrenia**

**Hypertension**

**Type 1 Diabetes**

**(N) Network alone** is sufficient to identify reference genes

**Crohn's Disease**

**Coronary Artery Disease**

**(S) Scores alone** are sufficient to identify reference genes

**Rheumatoid Arthritis**

**Bipolar Disorder**

**Type 2 Diabetes**

**(B) Both network and scores** help identify reference genes

Scores Only    Network Propagation    Network Only (PageRank)

# NetMix2 results on diseases where both network and scores help



Rheumatoid Arthritis | Bipolar Disorder | Type 2 Diabetes

**(B) Both network and scores** help identify reference genes

Scores Only · Network Propagation · NetMix2 · Network Only (PageRank)

NetMix2 outperforms network propagation on 2/3 diseases

# Anomaly detection



Normal means: Data $X_1, \ldots, X_n$ independently distributed as

$$X_i \sim \begin{cases} N(\mu, 1) & \text{if } i \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$
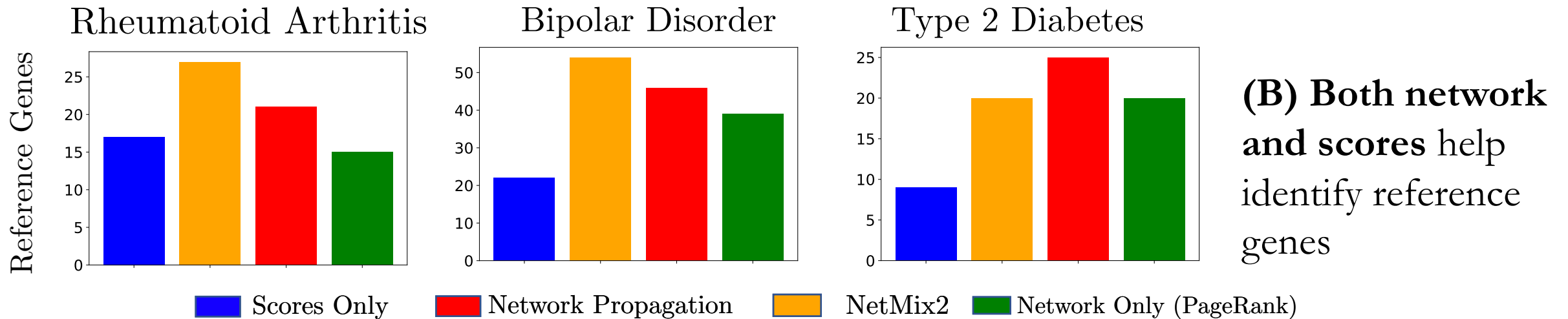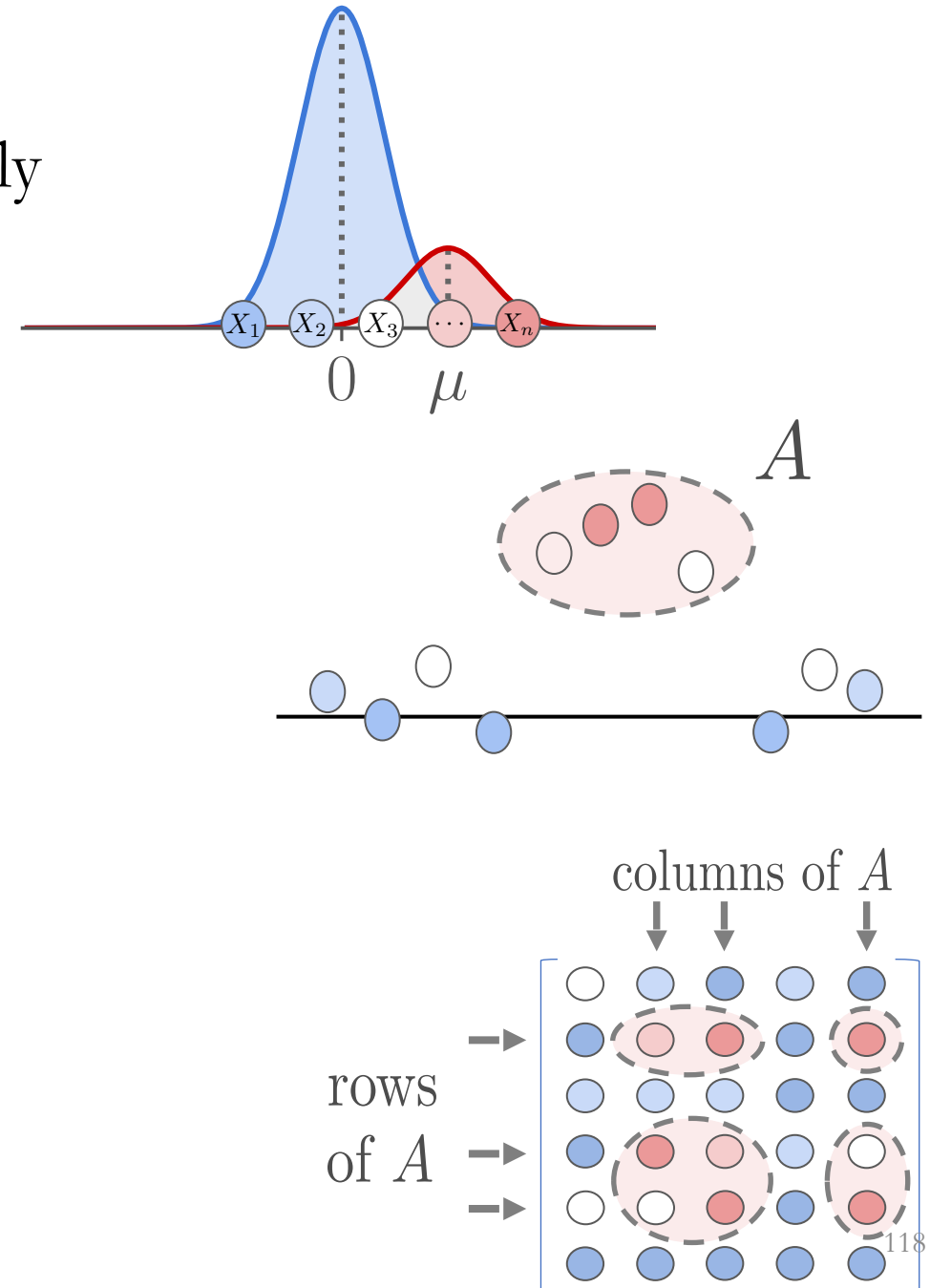
for anomaly A

Long history in statistics/ML:

- Unstructured anomalies: *Localfdr/empirical Bayes* methods (e.g. Efron et al., JASA 2001/2004, *Annals of Stats* 2007, etc), *Higher criticism* (Donoho and Jin, *Annals of Stats* 2004, etc), …

- Structured anomalies
  - **Intervals:** Jeng et al (JASA 2010)
  - **Submatrices:** Kolar et al (NeurIPS 2011), Chen and Xu (ICML 2014), Brennan et al (COLT 2018), Liu and A-C (KDD 2019)
  - **Connected subgraphs:** Qian et al (NeurIPS 2014), Aksoylar et al (ICML 2017), Cadena et al (AAAI 2018/TKDD 2019)
  - **Subgraphs w/ small cut**: Sharpnack et al (NeurIPS 2013/AISTATS 2013)
  - **Other:** Brennan et al (ICML 2020)
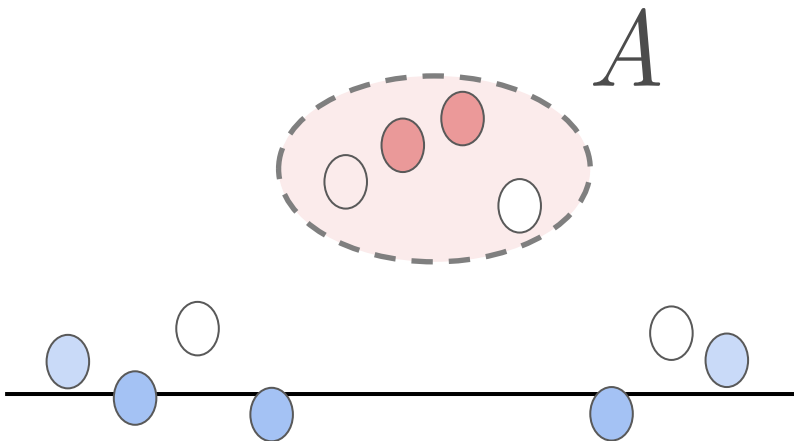
# Generalizing to anomaly detection

- $\mathcal{S}$ is **anomaly family** (set of subsets of $\{1, \dots, n\}$)
- $A \in \mathcal{S}$ is the anomaly

## Examples of anomaly families:

Interval family

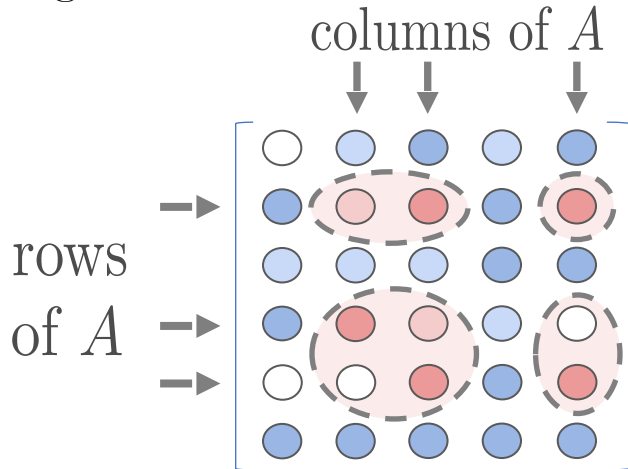$\mathcal{S} = \mathcal{I}_n$ = intervals $\{i, i+1, \dots, j\}$

*Changepoint detection*

Submatrix family
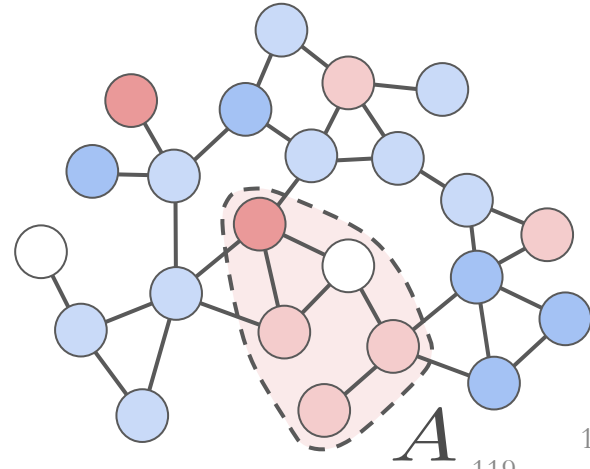
$\mathcal{S} = \mathcal{M}_N$ = submatrices of N

*Bi-clustering*

Connected family

$\mathcal{S} = \mathcal{C}_G$ = connected subgraphs of G

*Network anomaly detection*

# Anomalous Subset
~~Altered Subnetwork~~ Distribution

- n datapoints (e.g. vertices of interaction network)
- $\mathcal{S}$ is **anomaly family** (set of subsets of $\{1, \ldots, n\}$)
- $A \in \mathcal{S}$ is the anomaly

Datapoints $(X_1, \ldots, X_n)$ distributed as $X_i \sim \begin{cases} N(\mu, 1) & \text{if } i \in A \\ N(0, 1) & \text{otherwise} \end{cases}$

**Anomalous Subset Problem (ASP)**: Given data $(X_1, \ldots, X_n)$ and anomaly family $\mathcal{S}$, find anomaly A.

**Maximum Likelihood Estimator (MLE):**

$$\widehat{A}_{\text{MLE}} = \arg\max_{S \in \mathcal{S}} \frac{1}{\sqrt{|S|}} \sum_{i \in S} X_i$$

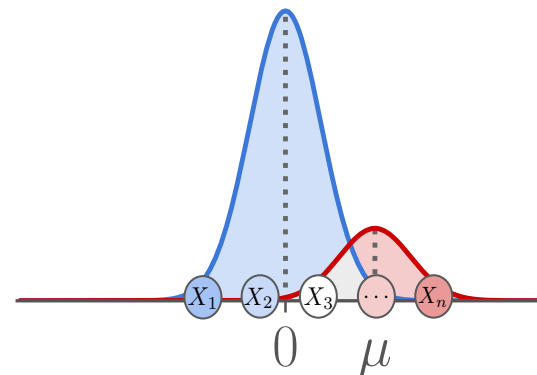# MLE is optimal for some anomaly families but not others

- Jeng et al (JASA 2010) show (asymptotic) "*near-optimality*" for <u>interval family</u> $\mathcal{S} = \mathcal{I}_n$

- Liu and A-C (KDD 2019) show similar guarantees for <u>submatrix family</u> $\mathcal{M}_N$

But we showed that MLE is a **biased** estimator for the <u>connected family</u> $\mathcal{S} = \mathcal{C}_G$

**Question**: for which anomaly families $\mathcal{S}$ is MLE biased?

**Maximum Likelihood Estimator (MLE):**

$$\widehat{A}_{\mathrm{MLE}} = \arg\max_{S \in \mathcal{S}} \frac{1}{\sqrt{|S|}} \sum_{i \in S} X_i$$



Data $X_1, \ldots, X_n$ distributed as

$$X_i \sim \begin{cases} N(\mu, 1) & \text{if } i \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

where anomaly $A \in \mathcal{S}$ is a member of anomaly family $\mathcal{S}$

# Our contribution

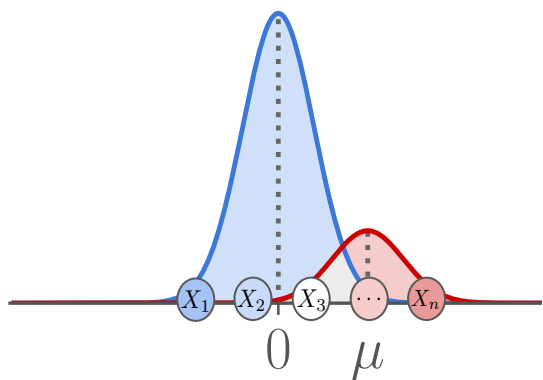Question: For which anomaly families $\mathcal{S}$ is the MLE biased?

**We show:** MLE is biased <--> number of sets in anomaly family $\mathcal{S}$ that contain the anomaly A is underlined exponential

Generalizes previous results on interval/submatrix family, which have sub-exponential size

Forward direction: ICML 2021 paper
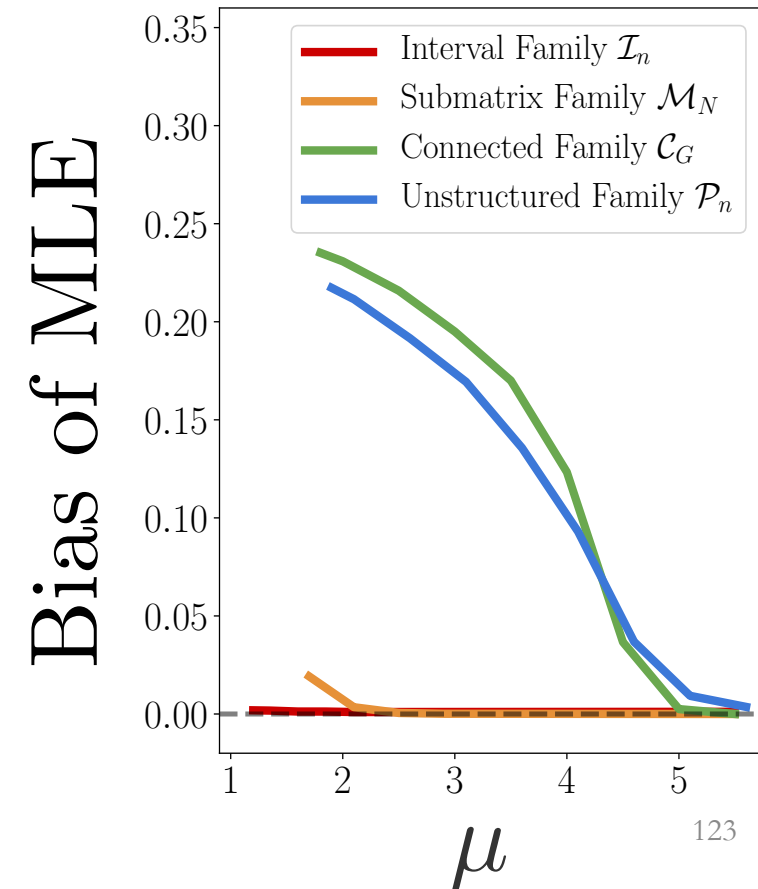Reverse direction: proved by Henri Schmidt+UC (unpublished)

Conjecture: result holds for exponential family distr. besides normal



Data $X_1, \ldots, X_n$ distributed as

$$X_i \sim \begin{cases} N(\mu, 1) & \text{if } i \in A \\ N(0, 1) & \text{otherwise} \end{cases}$$

where anomaly $A \in \mathcal{S}$ is a member of anomaly family $\mathcal{S}$



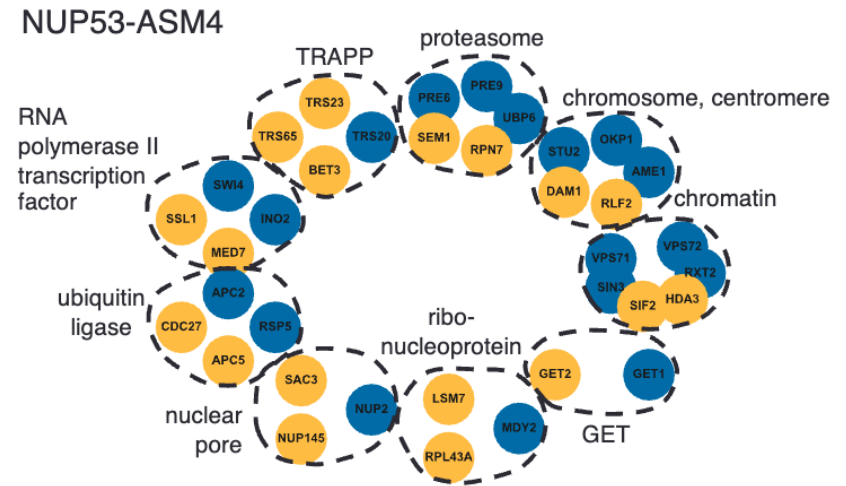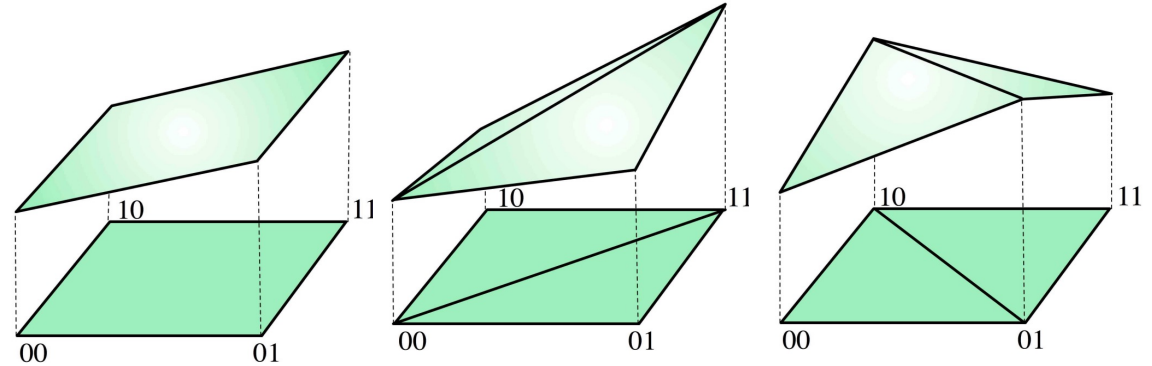Bias of MLE

- Interval Family $\mathcal{I}_n$
- Submatrix Family $\mathcal{M}_N$
- Connected Family $\mathcal{C}_G$
- Unstructured Family $\mathcal{P}_n$

# Learning genetic interactions (epistasis)



Brian Arnold

Ben Raphael

**Chitra\***, Arnold\*, Raphael. *In review at Nature Genetics.*
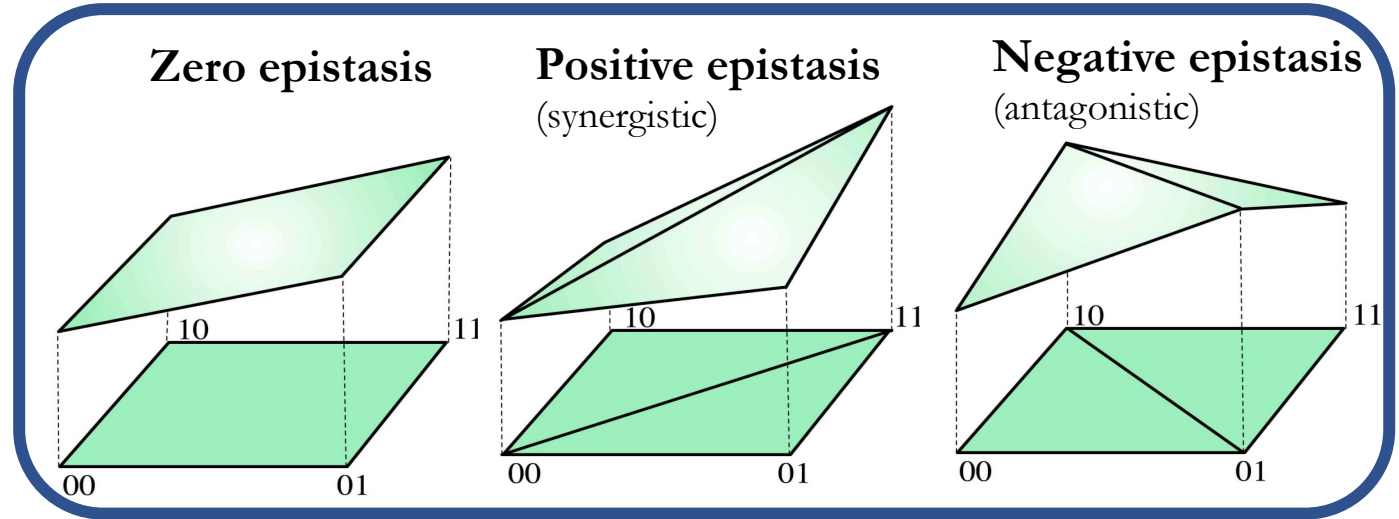
\* indicates joint first authorship

# Epistasis = genetic interactions - one gene mutation changes effect of other gene mutations

## Quantifying pairwise epistasis
(2 mutations)

Additive: $\epsilon = f_{11} - (f_{01} + f_{10})$

Multiplicative: $\epsilon = \dfrac{f_{11}}{f_{01} f_{10}}$



**Zero epistasis**

**Positive epistasis**
(synergistic)

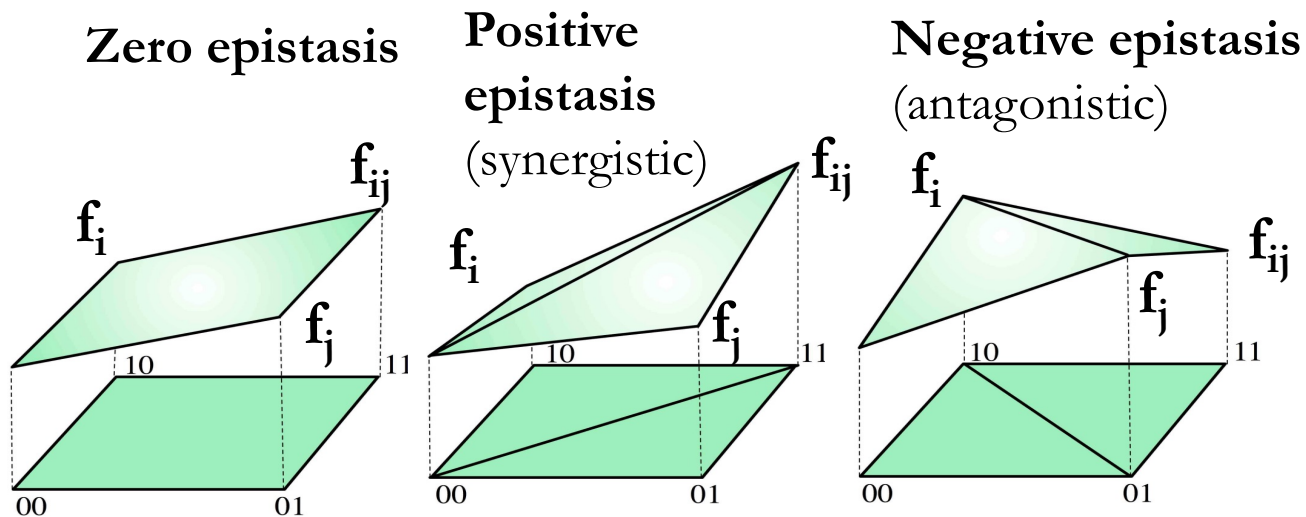**Negative epistasis**
(antagonistic)

# Epistasis = genetic interactions - one gene mutation changes effect of other gene mutations

Quantifying <u>pairwise</u> epistasis
(2 mutations)
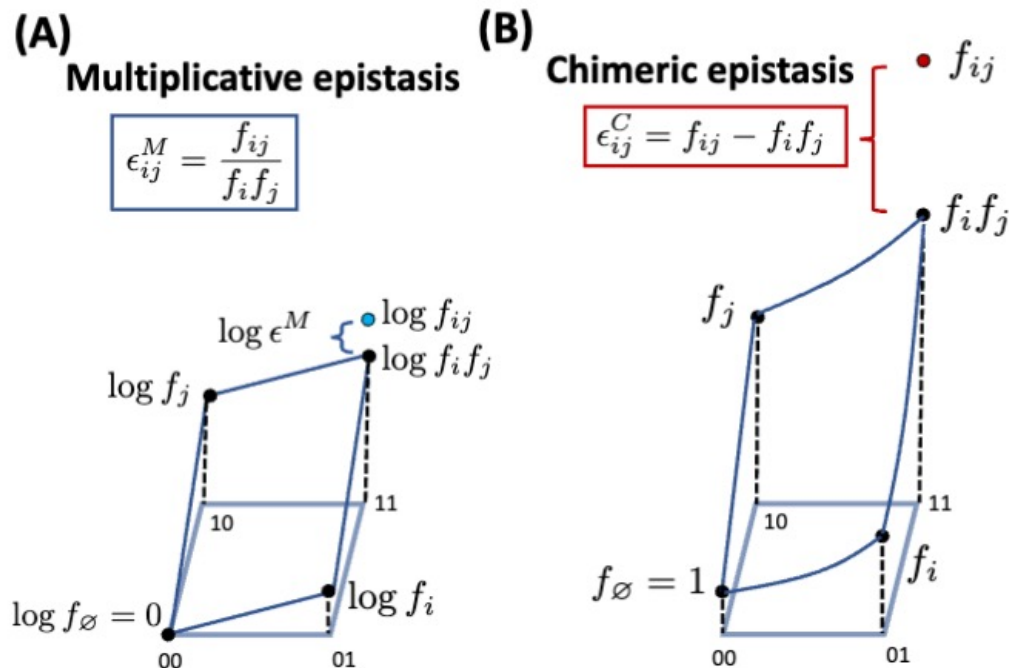
Additive: $\epsilon = f_{ij} - f_i - f_j$

Multiplicative: $\epsilon = \dfrac{f_{ij}}{f_i f_j}$

**Zero epistasis**

**Positive epistasis**
(synergistic)

**Negative epistasis**
(antagonistic)



Many papers in genetics claim to use multiplicative model but measure epistasis additively:
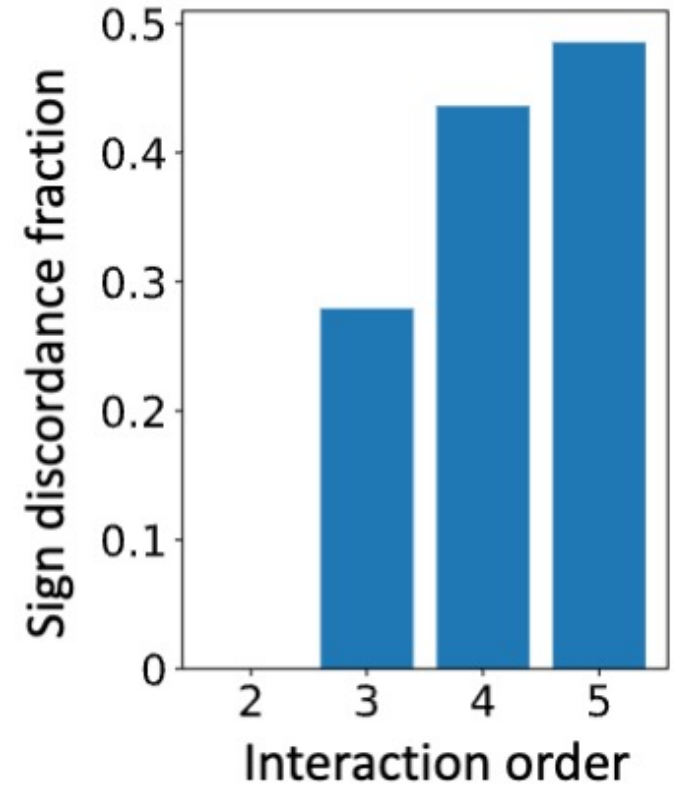
$$\epsilon^C = f_{ij} - f_i f_j$$

- *"Chimeric"* formula: a chimera of additive, multiplicative scales
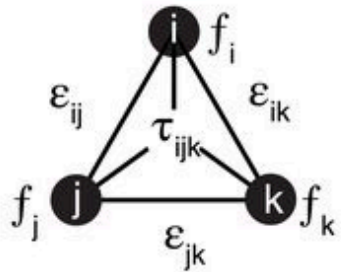- OK in practice: has same sign as multiplicative formula

**(A)** **Multiplicative epistasis**

$$\epsilon_{ij}^M = \frac{f_{ij}}{f_i f_j}$$

**(B)** **Chimeric epistasis**

$$\epsilon_{ij}^C = f_{ij} - f_i f_j$$

# Higher-order epistasis (3+ mutations)

Additive:
$$\epsilon_{ijk}^A = f_{ijk} - \left[ f_i + f_j + f_k + \epsilon_{ij}^A + \epsilon_{ik}^A + \epsilon_{jk}^A \right]$$
$$= f_{ijk} - f_{ij} - f_{ik} - f_{jk} + f_i + f_j + f_k.$$

Multiplicative:
$$\epsilon_{ijk}^M = \frac{f_{ijk}}{f_i f_j f_k \epsilon_{ij}^M \epsilon_{ik}^M \epsilon_{jk}^M} = \frac{f_{ijk} f_i f_j f_k}{f_{ij} f_{jk} f_{ik}}$$



Recent studies (*Science* 2018 + 2020) claim to use multiplicative fitness model but…

- Derive *"chimeric"* 3-way formula that combines additive, mult. Scales

- <u>No guarantees</u>: may have different sign versus multiplicative formula

**<span style="color:red">Hard to trust reported interactions!</span>**
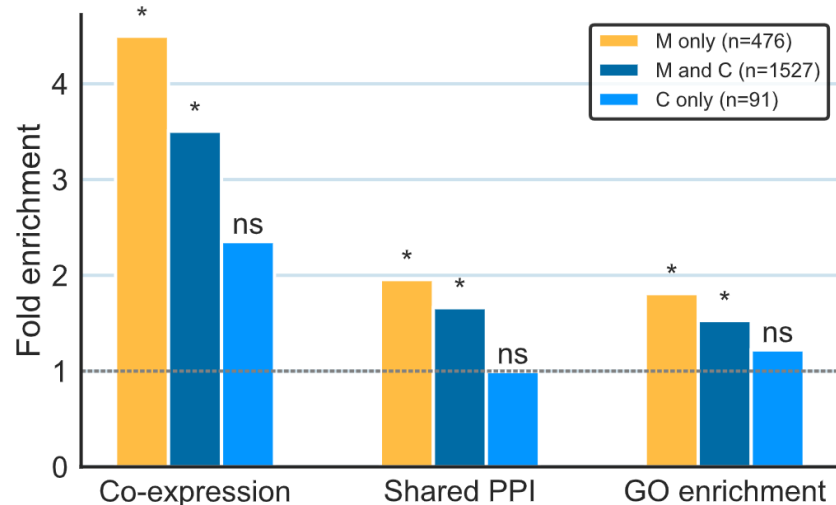
# Our contributions

1. <u>Unify different epistasis formulas using probabilistic framework</u>

   Epistasis formulas = different parametrizations of *multivariate Bernoulli* distribution (MVB)

   |  | **Fitness values f** | **Parameters of multivariate Bernoulli distribution** |
   |---|---|---|
   | **Additive epistasis measure $\epsilon^A$** | Log-probabilities $\log \mathbf{p}$ | Natural parameters $\boldsymbol{\beta}$ |
   | **Multiplicative epistasis measure $\epsilon^M$** | Probabilities $\mathbf{p}$ | Natural parameters $\boldsymbol{\beta}$ |
   | **Walsh coefficients** | Probabilities $\mathbf{p}$ | Moments of $(1 - 2X_1, \ldots, 1 - 2X_L)$ |
   | **Chimeric epistasis measure $\epsilon^C$** | Moments $\boldsymbol{\mu}$ | Joint cumulants $\kappa$ |

Our theory shows additive/multiplicative formulas are **more statistically sound** than chimeric formulas

2. Reanalyze *Science* data – learning 3-way interactions in yeast – using correct formula



Negative *(antagonistic)* 3-way interactions = functional redundancy
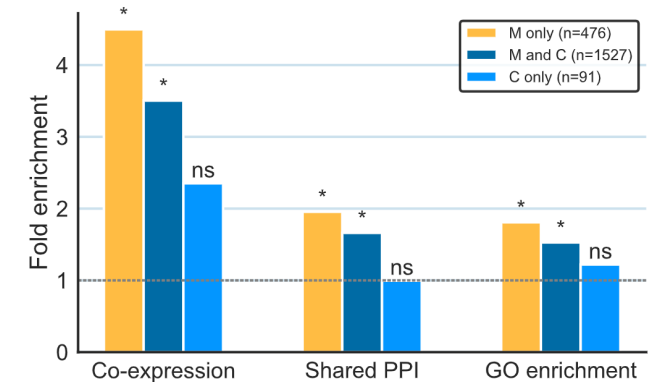
Using correct (mult.) formula finds ~500 more neg. interactions
- Significantly enriched for functional similarity measures
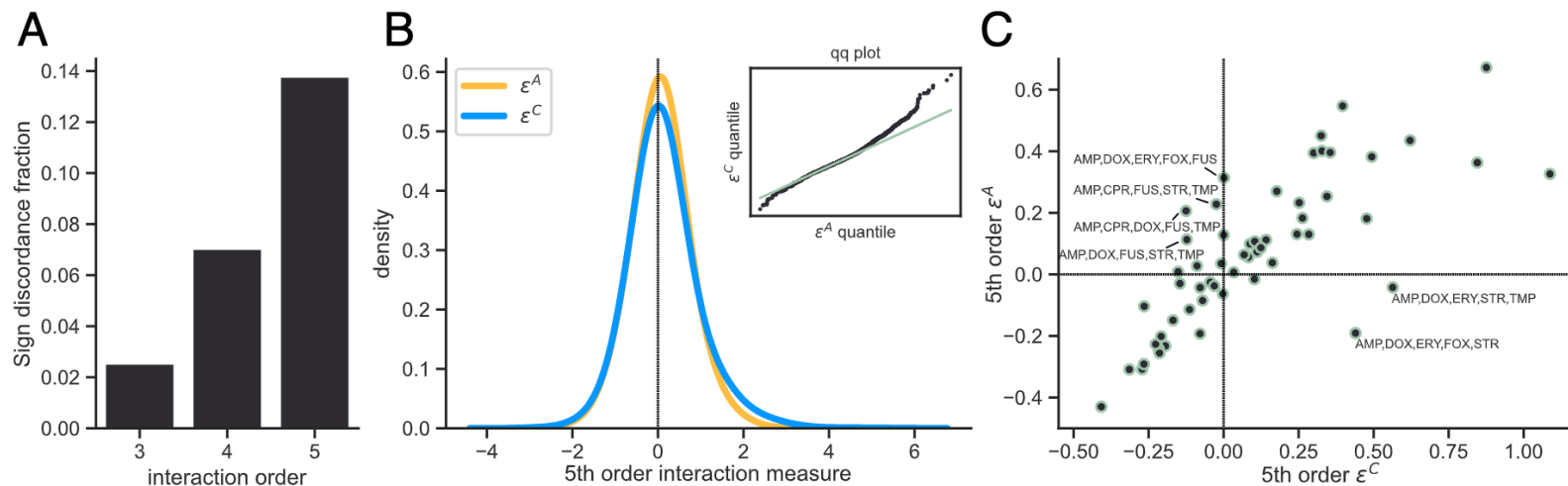- **extends trigenic interaction network by 25%**

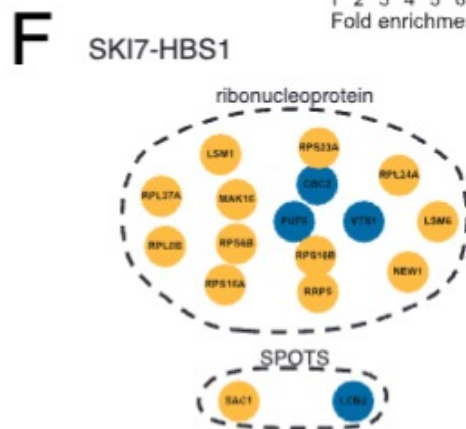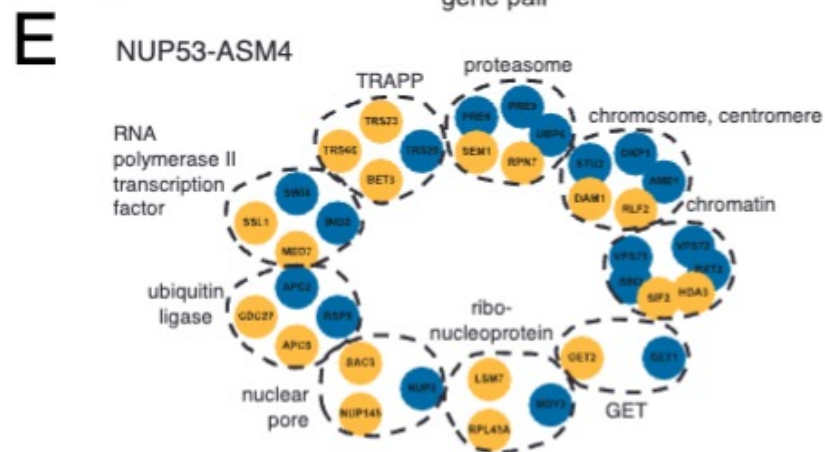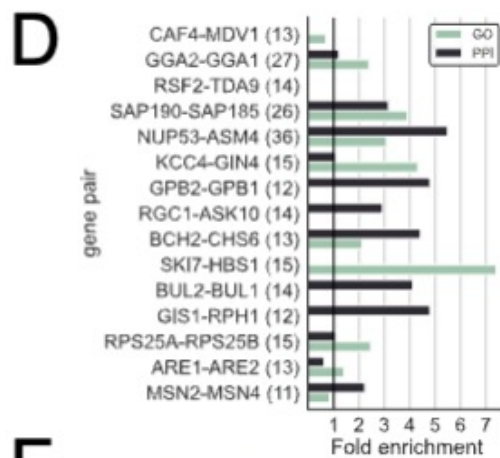# Sign disagreement leads to different biological findings

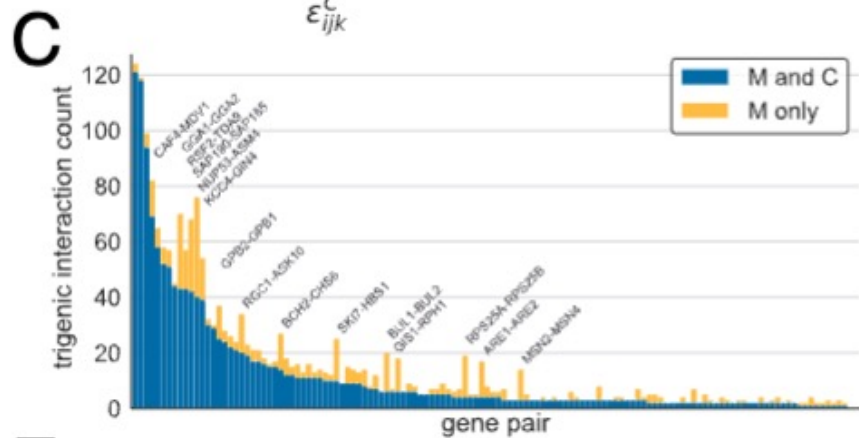**3-way epistasis in yeast**

|  |  | \multicolumn{3}{c}{Chimeric measure $\epsilon_{ijk}^{C}$} |  |  |
|---|---|---|---|---|
|  |  | Positive | Ambiguous | Negative |
| Multiplicative measure $\epsilon_{ijk}^{M}$ | Positive | 1197 | **259** | **0** |
|  | Ambiguous | **116** | 4291 | **91** |
|  | Negative | **10** | **466** | 1527 |



**Multi-way drug interactions**



131

Reanalysis of trigenic yeast interactions from Kuzmin et al. (*Science* 2018/2020)